

We thank all reviewers a lot for the comments and suggestions.

Reply to Referee #1

1. p9825: L14: One can almost answer the question "why does POP2 obtain different simulation results when changing the optimization level from O3 to anything else" by looking at the Intel compiler manual. The -O3 option uses a very fast math library which sacrifices accuracy for speed.

Response: Yes, it is true that compiler optimizations for speed may sacrifice accuracy. So it is reasonable that POP2 obtain different simulation results when changing the optimization level from O3. However, compiler optimizations should not bring new bugs to the applications. The example 2 in this manuscript shows that different results can not only result from accuracy sacrifices by the -O3 option but also result from compiler bugs that should be detected and avoided.

2. p9830, L10: Is this the same bug mentioned in Table 9? Please clarify.

Response: Yes, it is the same bug mentioned in Table 9. We modify the manuscript (P10 L15 - L16, P11 L9) to clarify.

3. p9838: In Table 3, please give a brief definition (or cite to the appropriate reference) of what the "fast|precise|strict|source" options do. Also please define SIMD.

Response: We modify the manuscript (P19) to complement these missing definitions.

4. pp9843-9845: If possible, please give brief description of why the bitwise-identical compiling setups obtained with Intel Fortran Compiler v11 are not reproducible with later Intel compilers.

Response: We add a brief description (P7 L6 - L12) about why the bitwise identical compiling setups significantly differ between Intel Compiler version 11 and later versions.

5. ... revisions to correct the many errors in English usage.

Response: We correct issues in English usage mentioned in the revised manuscript.

Reply to Referee #2

1. First, vectorization and SIMD are important options for most modern CPU architectures and they affect performance significantly for most ESMs, especially for OpenMP/MPI hybrid computation models. Therefore, the claim on P9832 L18 "such a compiler flag does not significantly decrease the computation performance" may be valid for some climate models only. But, it does not affect the validity of using the compiler flags from the biggest bitwise identical compiling setup set for code testing and code validating. In practice, it is always the best to select a compiling setup with the best computation performance while the scientific correctness of a model simulation is guaranteed.

Response: Thanks a lot for this comment. We agree with the importance of vectorization and SIMD. It is true that any selection of a compiler flag for a model simulation will not affect the code testing based on bitwise identical compiling setup sets. We improve the expression of the suggestion in the revised manuscript (P13 L12 – L24).

2. P9856 and P9857: the "Advantage of compiler flag B compared to compiler flag A" in both Figure 4 and Figure 5 are not well-defined. The axis labels on the right-hand sides of those figures are not mentioned anywhere.

Response: We complement it in the revised manuscript. (P39 L11 – L13, P40 L8 – L9)

3. P9843, P9844, P9845 and P9849: there is no explanation about why Intel compiler suite version 11 is significantly different with later versions in Tables 8, 9, 10 and 14.

Response: We add a brief description (P7 L6 - L12) about why the bitwise identical compiling setups significantly differ between Intel Compiler version 11 and later versions.

1 Bitwise Identical Compiling Setup: Prospective for 2 Reproducibility and Reliability of Earth System Modeling

3 R. Li^{2,1}, L. Liu^{1,3}, G. Yang^{2,1,3}, C. Zhang^{2,1} and B. Wang^{1,3,4}

4 [1]{Ministry of Education Key Laboratory for Earth System Modeling, Center for Earth System
5 Science (CESS), Tsinghua University, Beijing, China}

6 [2]{Department of Computer Science and Technology, Tsinghua University, Beijing, China}

7 [3]{Joint Center for Global Change Studies (JCGCS), Beijing, China}

8 [4]{State Key Laboratory of Numerical Modeling for Atmospheric Sciences and Geophysical
9 Fluid Dynamics (LASG), Institute of Atmospheric Physics, Chinese Academy of Sciences,
10 Beijing, China}

11 Correspondence to: L. Liu (liuli-cess@tsinghua.edu.cn), G. Yang (ygw@tsinghua.edu.cn)

12

13 Abstract

14 Reproducibility and reliability are fundamental principles of scientific research. A compiling
15 setup that includes a specific compiler version and compiler flags is an essential technical
16 supports for Earth system modeling. With the fast development of computer software and
17 hardware, a compiling setup has to be updated frequently, which challenges the reproducibility
18 and reliability of Earth system modeling. The existing results of a simulation using an original
19 compiling setup may be irreproducible by a newer compiling setup because trivial round-off
20 errors introduced by the change of compiling setup can potentially trigger significant changes
21 in simulation results. Regarding the reliability, a compiler with millions of lines of codes may
22 have bugs that are easily overlooked due to the uncertainties or unknowns in Earth system
23 modeling. To address these challenges, this study shows that different compiling setups can
24 achieve exactly the same (bitwise identical) results in Earth system modeling, and a set of
25 bitwise identical compiling setups of a model can be used across different compiler versions
26 and different compiler flags. As a result, the original results can be more easily reproduced; for
27 example, the original results with an older compiler version can be reproduced exactly with a
28 newer compiler version. Moreover, this study shows that new test cases can be generated based
29 on the differences of bitwise identical compiling setups between different models, which can

1 help detect software bugs ~~or risks~~ in the codes of models and compilers and finally improve the
2 reliability of Earth system modeling.

3 **1 Introduction**

4 Earth system modeling simulates interactions between components of the climate system (e.g.,
5 atmosphere, oceans, land surface, sea ice, etc.). It plays a critical role in understanding the past
6 and present climate, and in predicting future climate. An increasing number of models have
7 sprung up all over the world, including stand-alone component models and coupled models
8 consisting of multiple component models, such as Climate System Models (CSMs) and Earth
9 System Models (ESMs).

10 The development of models for Earth system modeling heavily depends on the advancement of
11 computer supports, not only in terms of hardware (such as high-performance computers) but
12 also in terms of software such as compiling setups that include compiler versions and compiler
13 flags. During the continuous ~~evolvment~~ evolution of the models, the compiling setups ~~for the~~
14 ~~model codes~~ have to be updated frequently for the compatibility ~~usage~~ of newer high-
15 performance computers with new processors and for better computing performance.

16 One may think it is easy to update compiling setups, ~~no more than~~ just by installing new
17 compiler version or changing compiler flags. However, it is challenging to update compiling
18 setups for Earth system modeling, because researchers may get significantly different results
19 from the same experiment when using different compiling setups (Liu et al., 2015b). A compiler
20 not only translates the codes in a high-level programming language to ~~the codes in~~ a low-level
21 language but also tries to improve computational ~~ing~~ performance of the codes with compiler
22 optimization schemes. Compilers from different families (for example, those in Table 1) and
23 different versions from the same compiler family ~~are generally different~~ generally differ in
24 performance optimization schemes as well as the corresponding implementations, ~~while~~ On
25 the other hand, different compiler flags of the same compiler version enables and disables
26 different sets of performance optimization schemes. That is why different compiling setups can
27 lead to different results of the same program. The updating of compiling setups therefore will
28 introduce at least two challenges to Earth system modeling. The first challenge ~~is about~~
29 ~~the concerns~~ reproducibility of simulation results. Due to the chaotic nature of the climate
30 system, more and more studies have shown that trivial round-off errors can trigger significant
31 changes in simulation results of Earth system modeling (Hong et al., 2013; Liu et al., 2015b;
32 Song et al., 2012). Due to the differences of performance optimization schemes ~~between~~ among

1 different compiling setups, a change of compiling setups potentially introduces round-off errors.
2 As a result, the ~~existing~~ results of a simulation ~~using an original~~ obtained with a compiling setup
3 may be irreproducible by another compiling setup.

4 The second challenge is ~~about~~ the reliability of the simulation results. Compilers are large-scale
5 programs with millions of lines of codes. It is well understood that with more lines of codes
6 there are more potential bugs in the program. Therefore, although there is generally a large
7 amount of software testing before releasing a compiler version, there ~~are still~~ can be unknown
8 bugs. Models for Earth system modeling are also large-scale numerical programs with ~~steadily~~
9 ~~increasing amount~~ more and more of code lines lines of code (Easterbrook and Johns, 2009).
10 There are already ESMs with nearly one million lines of codes (Alexander and Easterbrook,
11 2015). Therefore, it is possible that some bugs in a compiler version may be triggered by some
12 code segments in a model.

13 In response to these challenges, several issues about compiling setups should be
14 ~~concerned~~ investigated:

- 15 1) Can different compiling setups achieve the same (bitwise identical) simulation results? If
16 yes, it will be much easier to reproduce previous simulation results.
- 17 2) How can compiler flags to be selected ~~to select compiler flags when using a compiler~~ to
18 compile the codes of a model. Since a compiler version always contains many performance
19 optimization schemes, there are a lot of choices of compiler flags.
- 20 3) How to ~~find out~~ determine whether compiler bugs are triggered in a model simulation. If
21 compiler bugs can be detected, researchers can modify the code to avoid the compiler bugs
22 or select a “safer” compiling setup. Compiler bugs are very difficult to detect, especially
23 when they do not lead to a crash of the simulation. There are ~~a lot of~~ many uncertainties and
24 unknowns in Earth system modeling, so compiler bugs can easily be overlooked due to
25 these uncertainties or unknowns.

26 There are already efforts for the above-mentioned issues. It has been demonstrated that with a
27 certain compiler flag, different compiler versions can achieve bitwise identical simulation
28 results for a given model (Liu et al., 2015a). ~~—, while~~ But it is still not known whether ~~the~~
29 compiling setups with the same compiler version but different compiler flags can achieve
30 bitwise identical simulation results. In this paper, we call the compiling setups that can achieve
31 bitwise identical simulation results “bitwise identical compiling setups.” It is also unknown

1 whether the bitwise identical compiling setups of one model are appropriate for another model.
2 Baker et al. (2015) proposed a new ensemble-based consistency test for the Community Earth
3 System Model (CESM; Hurrell et al. 2013). It can effectively verify whether two compiling
4 setups can achieve consistent simulation results, especially when they do not achieve bitwise
5 identical simulation results. However, we cannot be sure whether a compiling setup is right or
6 wrong. In other words, it cannot help detect compiler bugs. As a result, it is possible that a
7 compiling setup with compiler bugs has been used for the development of a model for a number
8 of years, while a new compiling setup with bug fixes cannot be used for the model development
9 due to the failure in consistency tests.

10 The results in this paper show that the bitwise identical compiling setup sets of a model can ~~be~~
11 aerosextend to different compiler versions and different compiler flags. They can facilitate the
12 reproduction of original simulation results, ~~help assist~~ researchers to determine the compiler
13 flags for model simulations, help researchers build more test cases to detect bugs in models and
14 compilers, and finally improve the reproducibility and reliability of Earth system modeling.

15 The rest of this paper is organized as follows. Section 2 briefly introduces compiler
16 optimizations. Section 3 shows the bitwise identical compiling setups of three models. Section
17 4 uses examples to show what can be learned from the comparison of bitwise identical
18 compiling setups between different models. We conclude this paper with discussion in Section
19 5.

20 **2 Brief introduction to compiler optimizations**

21 Models for Earth system modeling are generally program~~m~~ed in languages such as Fortran, C
22 and C++. A number of compilers have been used for Earth system modeling, such as the
23 compiler families listed in Table 1. In the following context, we further introduce the Intel
24 compiler family and GNU Compiler Collection (GCC) with details.

25 The Intel compiler family, which is developed by the Intel Corporation, is a commercial
26 software product. It has been widely used for Earth system modeling, because most of the high-
27 performance computers for Earth system modeling are equipped with the CPUs manufactured
28 by the Intel Corporation. Table 2 shows the five latest Intel compiler versions (from version
29 11.1 released in 2009 to version 15.0.1 released in 2014). For each compiler version, there are
30 many compiler optimization options. Table 3 shows several compiler optimization options that
31 may impact the precision of floating-point calculation. They are common to all compiler

1 versions listed in Table 2. For a compiler flag such as “-fp-model”, there may be multiple
2 selections of the values.

3 GCC is the most widely used free compiler family in the world. Table 4 shows the five latest
4 GCC versions (from version 4.6.4 released in 2013 to version 5.1 released in 2015). For each
5 compiler version, there are also many compiler optimization options. Similar to Table 3, the
6 compiler optimization options in Table 5 may impact the precision of floating-point calculation
7 and are common to all GCC versions listed in Table 4.

8 **3 Bitwise identical compiling setups**

9 In this study, we use three models, namely, CAM5 (Neale et al., 2010), POP2 (Smith et al.,
10 2010) and FGOALS-g2 (Li et al., 2013a). To obtain bitwise identical compiling setups of a
11 given model, we should first design various compiling setups and then run the model using each
12 of them. In this section, we will briefly introduce the three models, the compiling setups and
13 the bitwise identical compiling setups of each model.

14 **3.1 Models and simulations**

15 The version of CAM5 used in this study is CAM5.3. It is released as the atmosphere component
16 of the CESM version 1.2 (CESM1.2). It contains more than 550,000 lines of source code mainly
17 programmed in Fortran. It can be used as a standalone model or the atmospheric component of
18 CESM1.2. In this study, we use CAM-5.3 as a standalone model. CAM-5.3 supports different
19 dynamic cores and different resolutions. To run the standalone CAM-5.3, we use the default
20 setting (details can be found at
21 http://www.cesm.ucar.edu/models/cesm1.2/cam/docs/ug5_3/ug.html), where the dynamic core
22 is finite volume and the resolution of the horizontal grid is 1.9°x2.5°.

23 POP2 used in this study is the ocean component of CESM-1.2. It is based on the POP version
24 2.1 of the Los Alamos National Laboratory. It contains more than 170,000 lines of source code
25 mainly programmed in Fortran. To run POP2 as a standalone model, we use the component set
26 “C_NORMAL_YEAR” of CESM-1.2, which uses POP2 as the ocean component and the other
27 components as data models. The horizontal grid selected is marked as “T62_gx1v6,” while the
28 other settings of the simulation are default.

29 FGOALS-g2 is a fully coupled CSM consisting of the atmosphere model GAMIL2 (Li et al.,
30 2013b), ocean model LICOM2 (Liu et al., 2004), land surface model CLM3 (Oleson et al.,

1 2004), and an improved version (Wang et al., 2009; Liu, 2010) of the sea ice model CICE4
2 (<http://climate.lanl.gov/Models/CICE>). It participated in the the Coupled Model
3 Intercomparison Project Phase 5 (CMIP5) and is widely used for scientific research. It contains
4 about 240,000 lines of source code mainly programmed in Fortran. GAMIL2 and CLM3 use
5 the same horizontal grid whose resolution is about 2.8°, while LICOM2 and CICE4 uses the
6 same horizontal grid whose resolution is about 1°. To run FGOALS-g2, we use the CMIP5 pre-
7 industry control (pi-Control) experiment setup.

8 All simulations of the models are run on the same high-performance computer named
9 Tansuo100 at Tsinghua University in China, which consists of more than 700 computing nodes,
10 each of which consists of two Intel Xeon 5670 6-core CPUs sharing 32GB main memory.
11 Specifically, we use 16, 16 and 17 processes to run CAM5.3, POP2 and FGOALS-g2,
12 respectively.

13 **3.2 Compiling setups**

14 ~~Through~~ By combining different settings of different compiler optimization options listed in
15 Table 3, there are more than 4,000 possible compiler flags. Considering there are four major
16 optimization levels (O0-O3) in an Intel compiler version, there are more than 16,000 possible
17 compiler flags for an Intel compiler version. Similarly, there are more than 1,000 possible
18 compiler flags for a GCC compiler version.

19 It is impractical for us to investigate all compiling setups. We decided to use five Intel compiler
20 versions (versions 11.1, 12.1, 13.0, 14.0.1, and 15.0.1) and five GCC compilers versions
21 (versions 4.6.4, 4.7.4, 4.8.5, 4.9.3, and 5.1) for this study, and take into consideration four
22 optimization levels (O0-O3). For a compiler version at an optimization level, we selected a
23 small number of compiler flags (Table 6 for the Intel compilers and Table 7 for the GCC
24 compilers).

25 **3.3 Bitwise identical compiling setups of models**

26 To obtain the bitwise identical compiling setups of a model (CAM5, POP2, or FGOALS-g2),
27 we use each compiling setup (in Section 3.2) to compile the model code and then run the
28 corresponding model simulation. A short integration is enough to check bitwise identity of
29 simulation results (Easterbrook and Johns, 2009). In detail, we use five model days for each
30 simulation and use the binary formatted data file of daily output of fields for bitwise identical

1 comparison. Tables 8-10 show the bitwise identical compiling setups of a model when using
2 the Intel compiler versions, while Tables 11-13 correspond to the GNU compiler versions. In
3 each table, the compiling setups corresponding to the same color (except for ~~the white color~~) of
4 simulation results constitute a bitwise identical compiling setup set of the same model. There
5 is no bitwise identical compiling setup set across the two compiler families.

6 [Given the Intel compiler versions, we can see that there is no bitwise identical compiling setup](https://software.intel.com/en-us/forums/intel-visual-fortran-compiler-for-windows/topic/281713)
7 [set between the version 11 and any other version. This is because the version 11 and the](https://software.intel.com/en-us/forums/intel-visual-fortran-compiler-for-windows/topic/281713)
8 [subsequent versions use different default instructions to generate the binary code](https://software.intel.com/en-us/forums/intel-visual-fortran-compiler-for-windows/topic/281713)
9 [\(https://software.intel.com/en-us/forums/intel-visual-fortran-compiler-for-](https://software.intel.com/en-us/forums/intel-visual-fortran-compiler-for-windows/topic/281713)
10 [windows/topic/281713\)](https://software.intel.com/en-us/forums/intel-visual-fortran-compiler-for-windows/topic/281713), which produces different bitwise results
11 [\(https://software.intel.com/en-us/forums/intel-visual-fortran-compiler-for-](https://software.intel.com/en-us/forums/intel-visual-fortran-compiler-for-windows/topic/279705)
12 [windows/topic/279705\)](https://software.intel.com/en-us/forums/intel-visual-fortran-compiler-for-windows/topic/279705).

13 **4 Comparison of bitwise identical compiling setup sets between models**

14 From Tables 8-10 (or Tables 11-13), we can find that, given the same compiler family, bitwise
15 identical compiling setup sets of different models are obviously different. What causes such
16 differences and what can we learn from the differences? To answer these questions, we take the
17 compiling setups of Intel compilers as an example. Based on the results in Tables 8-10, we can
18 generate ideal bitwise identical compiling setup sets (Table 14), following the criterion that if
19 any model achieves bitwise identical results with two different compiling setups, these
20 compiling setups belong to the same ideal bitwise identical compiling setup set. Through
21 comparing Tables 8-10 to Table 14, we can pose a number of questions; for example:

- 22 1) Regarding all Intel compiler versions, given compiler flag 2 (or 4), why does CAM5 obtain
23 different simulation results when changing compiler optimization level from O0 (or O1) to
24 O2 (or O3)?
- 25 2) Regarding Intel compiler version 13, why does POP2 obtain different simulation results
26 when changing the compiler optimization level from O3 to another level?
- 27 3) Regarding Intel compiler version 12, given optimization level O2, why does POP2 obtain
28 different simulation results when changing the compiler flag from 2 (or 3) to 1.
- 29 4) Regarding Intel compiler version 13, given optimization level O3, why does POP2 obtain
30 different simulation results when changing the compiler flag from 8 (or 9) to 1?

1 5) Regarding Intel compiler versions 13, 14 and 15, why does POP2 obtain the bitwise
2 identical results when changing the compiler flag from 1 to 2 (or 4), but CAM5 and
3 FGOALS-g2 do not?

4 Next, we search for answers to the first two questions, namely, what causes such differences
5 and what can we learn from the differences.

6 **4.1 Methodology**

7 If a code segment can trigger different compiler optimizations under different compiling setups,
8 it may lead to different results in different compiling setups. In the rest of this paper, we call
9 this kind of a code segment “a compilation-sensitive code segment” and call a code file with
10 compilation-sensitive code segments “a compilation-sensitive code file.” A model for Earth
11 system modeling generally contains a large number of code segments. To reveal why a model
12 does not achieve bitwise identical results in two different compiling setups, a straightforward
13 way is to find out all compilation-sensitive code segments for further analysis. Given two
14 compiling setups (denoted as C_A and C_B) that do not achieve bitwise identical results for a
15 simulation, we propose three stages for the detection of compilation-sensitive code segments:

16 1) Detect the compilation-sensitive code files. A model generally contains a number of source
17 code files. In the compiling process of a model, we can use C_A to compile a part-portion of
18 source code files while use C_B to compile the remaining source code files if the objective
19 files can be linked together. For example, at the first step, we can use C_A to compile all
20 source code files and then run a simulation to generate a reference result. At the second step,
21 we can divide the source code files into two parts, each of which takes about a half, and
22 then use different compiling setups to compile the two parts (use C_A to compile the first part
23 and use C_B to compile the second part, or use C_B to compile the first part and use C_A to
24 compile the second part). If the result from the same simulation is not bitwise identical with
25 the reference result, the part that is compiled with C_B should contain compilation-sensitive
26 code files, and next we will recursively detect compilation-sensitive code files in that part.

27 2) Detect compilation-sensitive code segments in a compilation-sensitive code file. We
28 propose to log (in binary format) and then bit-to-bit match the values of the input variables
29 and output variables of each code segment in the two compiling setups (C_A and C_B). A code
30 segment with bitwise identical inputs but different outputs is a compilation-sensitive code
31 segment. The size of a compilation-sensitive code segment should be as small as possible,

1 in order to facilitate further analysis. For a ~~code file~~source file that contains a large number
2 ~~of code lines~~containing many lines of code, we can either divide it into several new ~~code~~
3 files of smaller size and then repeat the first and second stages for these new files, or into
4 several big code segments at the first step and then recursively repeat the second stage for
5 the code segments that are compilation-sensitive. The size of a code segment cannot be too
6 small because the function calls for logging the values of variables may result in changes to
7 compiler optimizations so as to change simulation results. In other words, the splitting of a
8 code file or the inserting of the functions for logging values must keep bitwise simulation
9 results.

- 10 3) Analyze why a code segment is sensitive. In this stage, we should read the code to check
11 whether there are bugs. Sometimes, it is necessary to compare the differences of assembly
12 codes of the code segment under the two compiling setups.

13 Researchers may have to conduct the second and third stages manually. However, for the first
14 stage, we designed and implemented a software tool named CoSFID, which stands for
15 **C**ompilation-**S**ensitive code **F**ile **D**etection tool; it can automatically detect compilation-
16 sensitive code files (Section 4.2).

17 4.2 The CoSFID

18 Figure 1 shows the flowchart of CoSFID. The inputs include the two compiling setups (C_A and
19 C_B), the rules to compile and run the model, and the rules to compare results at bitwise identical
20 level. The outputs are a list of compilation-sensitive code files. CoSFID first generates a
21 reference result with the compiling setup C_A to compile all code files. Following the idea of the
22 first stage introduced in Section 4.1, CoSFID compiles and runs the model many times and
23 alternatively changes the compiling setup between C_A and C_B for some code files each time.

24 The biggest challenge to the design and implementation of CoSFID is how to control the
25 compilation process of each code file. A straightforward approach is to develop a common tool
26 that can successfully compile any model. However, this approach seems impractical because
27 different models may have different systems to compile the code, for example, using different
28 ways to specify code files and different ways to generate header files. We therefore propose to
29 use the original compiling system of a model and design a compiler wrapper accordingly. The
30 compiler wrapper is some script in CoSFID, which can replace the original compiler commands
31 used for compiling the model. For example, given that a model uses the Intel compiler

1 commands (i.e., `icc`, `icpc` and `ifort`) to compile the code, users should generate pseudo compiler
2 commands with the same names (i.e., `icc`, `icpc` and `ifort`) under a directory through symbolic
3 linking or copying the compiler wrapper of CoSFID, and then add the directory to the beginning
4 of the corresponding environment variable (for example, “PATH”) of the operating system to
5 make the pseudo compiler commands used for the compilation of the code, and then replace
6 the compiler flag for compiler optimizations by a label “-DCoSFID.” When compiling a code
7 file, CoSFID first gets the name of the file through the compiler wrapper; it then looks up the
8 current compiling setup for the file before switching the compiler version to the specified one
9 if necessary and using the specified compiler flag to replace the label “-DCoSFID”; it finally
10 compiles the code file.

11 4.3 Examples

12 4.3.1 Example 1

13 In this example, we search for the answer to the first question in Section 4: *regarding all Intel*
14 *compiler versions, given compiler flags 2 or 4, why does CAM5 obtain different simulation*
15 *results when changing compiler optimization level from O0 or O1 to O2 or O3? (as shown in*
16 *Table 8)* Following the methodology in Section 4.1, we first generate the two compiling setups
17 *C1* and *C2* using the Intel compiler version 13, compiler flag 2 (`-fp-model precise -fp-`
18 `speculation=strict -mp1 -no-vec -no-simd`) and two optimization levels (*O1* and *O2*); next, we
19 use CoSFID to find only one compilation-sensitive code file (`modal_aero_rename.F90`) from
20 more than 700 code files of CAM5. For further analysis, we split `modal_aero_rename.F90` into
21 two temporary code files, each of which contains only one subroutine, and then use CoSFID to
22 find that only the first subroutine (`modal_aero_rename_sub`) contains compilation-sensitive
23 code segments. Through logging and then comparing the values of input and output variables
24 of code segments in the two compiling setups, we find a compilation-sensitive code segment,
25 shown in Fig. 2. Given the same input (bitwise identical), this code segment can generate
26 slightly different results in different optimization levels (for example, Table 15). This is due to
27 the differences in assembly codes (Table 16). For the exponent *onethird* in Fig. 2, it is defined
28 as `1.0_r8/3.0_r8` in the program. The compiler optimization level *O1* will call function `pow` to
29 calculate the corresponding power function, while *O2* will intelligently find that the power
30 function is actually a cube root operation and then call `cbrrt` for the calculation.

1 After replacing variable *onethird* with $(1.0_r8/3.0_r8)$ throughout the code, CAM5 achieves
2 bitwise identical results with compiler flag 2 or 4 throughout all compiler optimization levels,
3 and finally the corresponding bitwise identical compiler setup sets of CAM5 are enlarged. For
4 example, the bitwise identical compiling setup set in green color and the set in blue color in
5 Table 8 are unified into one set.

6 4.3.2 Example 2

7 In this example, we search for the answer to the second question in Section 4: *regarding Intel*
8 *compiler version 13, why does POP2 obtain different simulation results when changing the*
9 *compiler optimization level from O3 to another level?* (as shown in Table 9) To generate the
10 two compiling setups *C1* and *C2*, we use the Intel compiler version 13, compiler flag 1 (*-fp-*
11 *model strict -fp-speculation=strict -mpl -no-vec -no-simd*) and two compiler optimization
12 levels (*O2* and *O3*). Using CoSFiD, we find only one compilation-sensitive code file
13 (*hmix_gm.F90*) from more than 500 code files of POP2. *hmix_gm.F90* contains about 10
14 subroutines and about 4,000 code lines. For further analysis, we split *hmix_gm.F90* into 10
15 temporary code files, each of which contains only one subroutine, and then use CoSFiD again
16 to find that only the temporary code file with the second subroutine (*hdiff_t_gm*) contains
17 compilation-sensitive code segments. Based on the binary values of input and output variables
18 of the code segments with the two compiling setups, we find a compilation-sensitive code
19 segment in the subroutine *hdiff_t_gm*, shown in Fig. 3. It is curious that given exactly the same
20 inputs, variable *WORK3* obtains significantly different results in the two compiling setups (for
21 example, Table 17). A manual result (Table 17) confirms correctness of the result in the
22 compiling setup with optimization level *O2*, but indicates that the code segment in Fig. 3
23 triggers a bug in the compiler when the compiler optimization level is *O3*.

24 It is almost impossible for us to fix a compiler bug. However, we can try to make the model
25 code not trigger the bug. Further analysis with assembly codes shows that the compiler performs
26 an optimization of loop fusion that merges four two-level loops at lines 1920-1999 of the code
27 file *hmix_gm.F90* into one loop. We intuitively guess that there are bugs in the loop fusion
28 optimization. To avoid the loop fusion optimization, we move the four two-level loops into a
29 new subroutine. Finally, POP2 achieves bitwise identical results with compiler flag 1
30 throughout all compiler optimization levels, and the corresponding bitwise identical compiling
31 setup sets of POP2 are enlarged. For example, the bitwise identical compiling setup set in red
32 and the set in green in Table 9 are unified into one set.

1 5 Discussion and conclusion

2 This study illustrates that a model can achieve bitwise identical results under different
3 compiling setups. For a given model, there are always a number of bitwise identical compiling
4 setup sets, some of which can be across not only different compiler flags but also different
5 versions of the same compiler family. As a result, the original results with an older compiler
6 version can be exactly reproduced with a newer compiler version. Moreover, the examples in
7 this paper reveal that bitwise identical compiling setup sets can be enlarged through carefully
8 modifying compilation-sensitive code segments, which will facilitate the exact reproduction of
9 original simulation results.

10 During the development of a model, the model codes increase continuously and need to be
11 tested frequently. The testing can be classified into two categories: scientific testing and
12 technical testing. Scientific testing, which is ~~through~~-evaluating the scientific meaning of
13 simulation results, is generally expensive, because it always requires long-~~time~~ simulations and
14 requires scientists to evaluate a large amount of ~~simulation~~ results. In contrast, technical testing,
15 which does not depend on the scientific meaning of simulation results, is generally cheap. For
16 example, short simulations (such as several model days) are enough for bitwise identical testing,
17 and bitwise identical testing can be conducted automatically without any burden to scientists
18 (Easterbrook and Johns, 2009). Technical testing therefore should be much more frequent than
19 scientific testing. Since a bitwise identical compiling setup set contains a number of compiling
20 setups that should achieve exactly the same results for a model simulation, it can bring more
21 cases for technical testing. For example, given that a new code version evolves from an old
22 code version with new modifications, the bitwise identical compiling setup sets of each code
23 version can be obtained automatically. If the two code versions do not have the same bitwise
24 identical compiling setup sets, new test cases can be generated for checking why this happens,
25 for example because of bugs in the codes or compilation-sensitive code segments. If there are
26 compilation-sensitive code segments in the new modifications, we advise researchers to make
27 them insensitive, to make each bitwise identical compiling setup set as big as possible for
28 further development of the model. The first example in Section 4.3 reveals that a compilation-
29 sensitive code segment can become insensitive after a slight code modification.

30 Although the bitwise identical compiling setup sets of different models are generally different,
31 the differences can effectively bring more test cases to detect software bugs in model
32 simulations, especially the bugs of compilers. Although scientists of Earth system modeling

1 generally cannot modify the code of a compiler to fix a bug, they can modify the code of a
2 model to make sure that the model code will not trigger a compiler bug again. For example,
3 based on the differences of bitwise identical compiling setup sets among different models
4 (CAM5, POP2 and FGOALS-g2), we found that a code segment of POP2 triggers a bug of the
5 Intel compiler version 13, and the compiler bug will not be triggered again with a slight
6 modification to the code segment.

7 There are generally a large number of choices of compiler flags. Researchers may tend to select
8 a compiler flag that can achieve the best computation performance for a model simulation. Our
9 performance evaluation shows that the compiler flag 3 can achieve the best computation
10 performance among the compiler flags in Table 6. According to Tables 8-10, the bitwise
11 identical compiling setup set corresponding to compiler flag 3 is small. It is already known that
12 climate simulation results can be sensitive to round-off errors. ~~For the simulations that are~~
13 ~~sensitive to round-off errors, the simulation results are either irreproducible or bitwise~~
14 ~~identically reproducible (Liu et al., 2015b).~~ To make simulation results most easily reproduced,
15 ~~we suggest~~ researchers may be able to use the compiler flag of the best computation
16 performance in ~~thea~~ biggest bitwise identical compiling setup set for a model simulation, when
17 the change of compiler flags will not significantly decrease the computation performance. For
18 example, ~~we propose to~~ researchers can use the compiler flag “*-O3 -fp-model strict -fp-*
19 *speculation=strict -mp1 -no-vec -no-simd*” for the simulation of the atmosphere models CAM5
20 and GAMIL2-model when the Intel compilers are used, ~~because. As shown in Fig. 4 and Fig.~~
21 ~~5,~~ such a compiler flag does not significantly decrease the computation performance, especially
22 when the number of processes is big (~~shown in Fig. 4 and Fig. 5~~). Please note that any selection
23 of a compiler flag for a model simulation will not affect the code testing based on bitwise
24 identical compiling setup sets.

25

26 **Code availability**

- 27 1. The source code of CESM version 1.2 can be obtained at
28 <http://www.cesm.ucar.edu/models/cesm1.2/>.
- 29 2. The source code of FGOALS-g2 is currently not publicly available. You can contract us for
30 more information.
- 31 3. The source code of CoSFiD is available at <https://github.com/liruizhe/CoSFiD>.

1 4. The compilation-sensitive code files mentioned in Section 4.3 will be included in the
2 supplement.

3 **Acknowledgements**

4 This work is supported in part by the Natural Science Foundation of China (no. 41275098),
5 the National Grand Fundamental Research 973 Program of China (no. 2014CB441302) and
6 the Tsinghua University Initiative Scientific Research Program (no. 20131089356).

7

1 **References**

- 2 Alexander, K., & Easterbrook, S. M.: The software architecture of climate models: a
3 graphical comparison of CMIP5 and EMICAR5 configurations. *Geosci. Model Dev.*, 8, 1221-
4 1232, 2015.
- 5 Baker, A. H., Hammerling, D. M., Levy, M. N., Xu, H., Dennis, J. M., Eaton, B. E., ... &
6 Williamson, D.: A new ensemble-based consistency test for the Community Earth System
7 Model. *Geosci. Model Dev. Discuss.*, 8, 3823-3859, 2015.
- 8 Easterbrook, S. M. and Johns, T. C.: Engineering the software for understanding climate change,
9 *Comput. Sci. Eng.*, 11, 65–74, 2009.
- 10 Hong, S. Y., Koo, M. S., Jang, J., Esther Kim, J. E., Park, H., Joh, M. S., ... & Oh, T. J.: An
11 Evaluation of the Software System Dependency of a Global Atmospheric model, *Mon. Weather*
12 *Rev.*, 141, 4165–4172, 2013.
- 13 Hurrell, J. W., Holland, M. M., Gent, P. R., Ghan, S., Kay, J. E., Kushner, P. J., ... & Marshall,
14 S.: The community earth system model: a framework for collaborative research. *Bulletin of the*
15 *American Meteorological Society*, 94(9), 1339-1360, 2013.
- 16 Li, L., Lin, P., Yu, Y., Wang, B., Zhou, T., Liu, L., ... & Qiao, F.: The flexible global ocean-
17 atmosphere-land system model, Grid-point Version 2: FGOALS-g2. *Advances in Atmospheric*
18 *Sciences*, 30, 543-560, 2013a.
- 19 Li, L., Wang, B., Dong, L., Liu, L., Shen, S., Hu, N., ... & Yang, G.: Evaluation of grid-point
20 atmospheric model of IAP LASG version 2 (GAMIL2). *Advances in Atmospheric Sciences*,
21 30, 855-867, 2013b.
- 22 Liu, H. L., Zhang, X.H., Li, W., Yu, Y.Q., Yu, R.C.: A eddy-permitting oceanic general
23 circulation model and its preliminary evaluations, *Adv. Atmos. Sci.* 21, 675–690, 2004.
- 24 Liu, L., Li, R., Zhang, C., Yang, G., Wang, B., and Dong, L.: Enhancement for bitwise
25 identical reproducibility of Earth system modeling on the C-Coupler platform, *Geosci. Model*
26 *Dev. Discuss.*, 8, 2403-2435, doi:10.5194/gmdd-8-2403-2015, 2015a.
- 27 Liu, L., Peng, S., Zhang, C., Li, R., Wang, B., Sun, C., ... & Yang, G.: Importance of bitwise
28 identical reproducibility in earth system modeling and status report, *Geosci. Model Dev.*
29 *Discuss.*, 8, 4375–4400, 2015b.

- 1 Liu, J.: Sensitivity of sea ice and ocean simulations to sea ice salinity in a coupled global climate
2 model. *Science China Earth Sciences*, 53(6), 911-918, 2010.
- 3 Neale, R. B., Chen, C. C., Gettelman, A., Lauritzen, P. H., Park, S., Williamson, D. L., ... &
4 Taylor, M. A.: Description of the NCAR community atmosphere model (CAM 5.0). NCAR
5 Tech. Note NCAR/TN-486+ STR, 2010.
- 6 Oleson, K. W., Dai, Y., Bonan, G., Bosilovich, M., Dickinson, R., Dirmeyer, P., ... & Zeng, X.:
7 Technical description of the community land model (CLM). NCAR Technical Note NCAR/TN-
8 461+ STR, National Center for Atmospheric Research, Boulder, CO, 2004.
- 9 Smith, R., Jones, P., Briegleb, B., Bryan, F., Danabasoglu, G., Dennis, J., ... & Yeager, S.: The
10 Parallel Ocean Program (POP) Reference Manual Ocean Component of the Community
11 Climate System Model (CCSM) and Community Earth System Model (CESM). Rep. LAUR-
12 01853, 141, 2010.
- 13 Song, Z., Qiao, F., Lei, X., and Wang, C.: Influence of parallel computational uncertainty on
14 simulations of the Coupled General Climate Model, *Geosci. Model Dev.*, 5, 313–319,
15 doi:10.5194/gmd-5-313-2012, 2012.
- 16 Wang, X. C., Liu, J. P., Yu, Y. Q., Liu, H. L., & Li, L. J.: Numerical simulation of polar climate
17 with FGOALS-g1. 1. *Acta Meteorologica Sinica*, 67(6), 961-972, 2009.

1 Table 1. Compiler families used for Earth system modeling. They are from the supported
 2 compiler lists of several ESMs.

Compiler family	Free or commercial	Supported hardware platforms	Supported programming languages
GNU	Free	Almost all common platforms	Fortran, C, C++, etc.
Intel	Commercial	x86 and x86-64bit architectures	Fortran, C, C++
PGI	Commercial	x86, x86-64bit, CUDA, and ARM architectures	Fortran, C, C++
Lahey	Commercial	x86 and x86-64bit architectures	Fortran
PathScale EKOPath	Commercial	x86 and x86-64bit architectures	Fortran, C, C++
Cray	Commercial	Cray supercomputer series (x86, x86-64bit and CUDA architectures)	Fortran, C, C++

3

1 Table 2. Five latest versions of the Intel compilers.

Compiler version	Release date
11.1	June 23, 2009
12.1	September 8, 2011
13.0	September 5, 2012
14.0.1	October 18, 2013
15.0.1	October 30, 2014

2

- 1 Table 3. Intel compiler optimization options that may impact the precision of floating-point
- 2 calculation. They are common to the compiler versions listed in Table 2.

Compiler optimization option	Description
-fp-model [fast precise strict] [source]	Controls the semantics of floating-point calculations: <u>fast: Enables more aggressive optimizations on floating-point data.</u> <u>precise: Enables value-safe optimizations on floating-point data.</u> <u>strict: Enables precise and except, disables contractions, and enables pragma stdc fenv_access.</u> <u>Source: Rounds intermediate results to source-defined precision and enables value-safe optimizations.-</u>
-fp-speculation fast precise safe strict	Tells the compiler the mode in which to speculate on floating-point operations. <u>fast: Tells the compiler to speculate on floating-point operations.</u> <u>safe: Tells the compiler to disable speculation if there is a possibility that the speculation may cause a floating-point exception.</u> <u>strict: Tells the compiler to disable speculation on floating-point operations.</u>
-mp1	Improves floating-point precision and consistency.
-[no-]vec	Enables or disables vectorization.
-[no-]simd	Enables or disables the SIMD <u>(Single instruction, multiple data)</u> vectorization feature of the compiler.
-[no-]fp-port	Rounds floating-point results after floating-point operations.

-[no-]ftz	Flushes denormal results to zero.
-pc[n]	Enables control of floating-point significant precision.
-[no-]prec-div	Improves precision of floating-point divides.
-[no-]prec-sqrt	Improves precision of square root implementations.

1

- 1 Table 4. Five latest versions of the GCC compilers. The release date of a given compiler version
2 in the table is the release date of its latest revision version.

Compiler version	Release date
4.6.4	April 12, 2013
4.7.4	April 13, 2013
4.8.5	June 23, 2015
4.9.3	June 26, 2015
5.1	April 22, 2015

3

- 1 Table 5. GCC compiler optimization options that may impact the precision of floating-point
 2 calculation. They are common to the compiler versions listed in Table 4.

Compiler flag	Description
-ffloat-store	Do not store floating-point variables in registers, and inhibit other options that might change whether a floating-point value is taken from a register or memory.
-ffast-math	Sets -fno-math-errno, -funsafe-math-optimizations, -ffinite-math-only, -fno-rounding-math, -fno-signaling-nans and -fcx-limited-range.
-f[no-]unsafe-math-optimizations	Allow optimizations for floating-point arithmetic that (a) assume that arguments and results are valid and (b) may violate IEEE or ANSI standards. When used at link-time, it may include libraries or startup files that change the default FPU control word or other similar optimizations.
-f[no-]associative-math	Allow re-association of operands in series of floating-point operations.
-f[no-]reciprocal-math	Allow the reciprocal of a value to be used instead of dividing by the value if this enables optimizations.
-f[no-]finite-math-only	Allow optimizations for floating-point arithmetic that assume that arguments and results are not NaNs or \pm Infs.
-f[no-]rounding-math	Disable transformations and optimizations that assume default floating-point rounding behavior.
-f[no-]cx-limited-range	When enabled, this option states that a range reduction step is not needed when performing complex division. Also, there is no checking whether the result of a complex multiplication or division is "NaN + I*NaN", with an attempt to rescue the situation in that case.

1 Table 6. Intel compiler flags that are based on the compiler optimization options given in Table
 2 3. The first compiler flag is the strictest one (which limits compiler optimizations most
 3 significantly), while every other compiler flag is derived from the first one through changing
 4 only one compiler optimization option.

No.	Compiler flag
1	-fp-model strict -fp-speculation=strict -mp1 -no-vec -no-simd
2	-fp-model precise -fp-speculation=strict -mp1 -no-vec -no-simd
3	-fp-model fast -fp-speculation=strict -mp1 -no-vec -no-simd
4	-fp-model source -fp-speculation=strict -mp1 -no-vec -no-simd
5	-fp-model strict -fp-speculation=safe -mp1 -no-vec -no-simd
6	-fp-model strict -fp-speculation=fast -mp1 -no-vec -no-simd
7	-fp-model strict -fp-speculation=strict -no-vec -no-simd
8	-fp-model strict -fp-speculation=strict -mp1 -vec -no-simd
9	-fp-model strict -fp-speculation=strict -mp1 -no-vec -simd

5

1 Table 7. GCC compiler flags that are based on the compiler optimization options listed in Table
 2 5. The first compiler flag is the strictest one (which limits compiler optimizations most
 3 significantly), while every other compiler flag is derived from the first one through changing
 4 only one compiler optimization option.

No.	Compiler flag
1	-ffloat-store -fno-unsafe-math-optimizations -fno-associative-math -fno-reciprocal-math -fno-finite-math-only -fno-rounding-math -fno-cx-limited-range
2	-fno-unsafe-math-optimizations -fno-associative-math -fno-reciprocal-math -fno-finite-math-only -fno-rounding-math -fno-cx-limited-range
3	-ffloat-store -funsafe-math-optimizations -fno-associative-math -fno-reciprocal-math -fno-finite-math-only -fno-rounding-math -fno-cx-limited-range
4	-ffloat-store -fno-unsafe-math-optimizations -fassociative-math -fno-reciprocal-math -fno-finite-math-only -fno-rounding-math -fno-cx-limited-range
5	-ffloat-store -fno-unsafe-math-optimizations -fno-associative-math -freciprocal-math -fno-finite-math-only -fno-rounding-math -fno-cx-limited-range
6	-ffloat-store -fno-unsafe-math-optimizations -fno-associative-math -fno-reciprocal-math -ffinite-math-only -fno-rounding-math -fno-cx-limited-range
7	-ffloat-store -fno-unsafe-math-optimizations -fno-associative-math -fno-reciprocal-math -fno-finite-math-only -frounding-math -fno-cx-limited-range
8	-ffloat-store -fno-unsafe-math-optimizations -fno-associative-math -fno-reciprocal-math -fno-finite-math-only -fno-rounding-math -fcx-limited-range

5

1 Table 8: Simulation results of CAM5 with various compiling setups of Intel compilers. The
 2 compiler flags are given in Table 6. Each color represents a bitwise identical result except for
 3 the white. A simulation result that emerges only once is in white color with a unique number.

Compiler optimization level	No. of compiler flag	Version of Intel compiler				
		11	12	13	14	15
-O0	1	Yellow	Green	Red	Red	Red
	2	Yellow	Green	Green	Green	Green
	3	(1)	Green	Green	Green	Green
	4	Yellow	Green	Green	Green	Green
	5	Yellow	Green	Red	Red	Red
	6	Yellow	Green	Red	Red	Red
	7	Yellow	Green	Red	Red	Red
	8	Yellow	Green	Red	Red	Red
	9	Yellow	Green	Red	Red	Red
-O1	1	Yellow	Green	Red	Red	Red
	2	Yellow	Green	Green	Green	Green
	3	(2)	(3)	(4)	(5)	(6)
	4	Yellow	Green	Green	Green	Green
	5	Yellow	Green	Red	Red	Red
	6	Yellow	Green	Red	Red	Red
	7	Yellow	Green	Red	Red	Red
	8	Yellow	Green	Red	Red	Red
	9	Yellow	Green	Red	Red	Red
-O2	1	Yellow	Green	Red	Red	Red
	2	Yellow	Blue	Blue	Blue	Blue
	3	(7)	(8)	(9)	(10)	(11)
	4	Yellow	Blue	Blue	Blue	Blue
	5	Yellow	Green	Red	Red	Red
	6	Yellow	Green	Red	Red	Red
	7	Yellow	Green	Red	Red	Red
	8	Yellow	Green	Red	Red	Red
	9	Yellow	Green	Red	Red	Red
-O3	1	Yellow	Green	Red	Red	Red
	2	Yellow	Blue	Blue	Blue	Blue
	3	(12)	(13)	(14)	(15)	(16)
	4	Yellow	Blue	Blue	Blue	Blue
	5	Yellow	Green	Red	Red	Red
	6	Yellow	Green	Red	Red	Red
	7	Yellow	Green	Red	Red	Red
	8	Yellow	Green	Red	Red	Red
	9	Yellow	Green	Red	Red	Red

4

1 Table 9: Similar to Table 8 except for the simulation results of POP2. Each table cell with "--
 2 "-- means that the compilation of POP2 fails under the corresponding compiling setup, due to
 3 issue DPD200178252 of Intel compilers ([https://software.intel.com/en-us/articles/intel-
 4 composer-xe-2013-compilers-fixes-list](https://software.intel.com/en-us/articles/intel-composer-xe-2013-compilers-fixes-list)).

Compiler optimization level	No. of compiler flag	Version of Intel compiler				
		11	12	13	14	15
-O0	1		--			
	2					
	3	(1)				
	4					
	5		--			
	6		--			
	7		--			
	8		--			
	9		--			
-O1	1		--			
	2					
	3	(2)			(3)	(4)
	4					
	5		--			
	6		--			
	7		--			
	8		--			
	9		--			
-O2	1		--			
	2					
	3	(5)			(6)	(7)
	4					
	5		--			
	6		--			
	7		--			
	8		--			
	9		--	(8)		
-O3	1		--			
	2					
	3	(9)	(10)	(11)	(12)	(13)
	4					
	5		--			
	6		--			
	7		--			
	8		--	(14)		
	9		--	(15)		

5

1 Table 10: Similar to Table 8 except for the simulation results of FGOALS-g2.

Compiler optimization level	No. of compiler flag	Version of Intel compiler				
		11	12	13	14	15
-O0	1					
	2					
	3	(1)				
	4					
	5					
	6					
	7					
	8					
	9					
-O1	1					
	2					
	3	(2)	(3)	(4)	(5)	(6)
	4					
	5					
	6					
	7					
	8					
	9					
-O2	1					
	2					
	3	(7)	(8)	(9)	(10)	(11)
	4					
	5					
	6					
	7					
	8					
	9					
-O3	1					
	2					
	3	(12)	(13)	(14)	(15)	(16)
	4					
	5					
	6					
	7					
	8					
	9					

2

1 Table 11: Simulation results of CAM5 with various compiling setups of GCC compilers. The
 2 compiler flags are given in Table 7. Each color represents a bitwise identical result except the
 3 white. A simulation result that emerges only once is in white color with a unique number.

Compiler optimization level	No. of compiler flag	Version of GCC compiler				
		4 6 4	4 7 4	4 8 5	4 9 3	5 1 0
-O0	1					
	2					
	3	(1)				
	4					
	5	(2)				
	6					
	7					
	8	(3)				
-O1	1					
	2					
	3	(4)	(5)			(6)
	4					
	5	(7)	(8)			(9)
	6					
	7					
	8					
-O2	1					
	2					
	3	(10)	(11)			(12)
	4					
	5	(13)	(14)	(15)	(16)	(17)
	6					
	7					
	8					
-O3	1					
	2					
	3	(18)	(19)			(20)
	4					
	5	(21)	(22)	(23)	(24)	(25)
	6					
	7					
	8		(26)			

4

1 Table 12: Similar to Table 11 except for the simulation results of POP2.

Compiler optimization level	No. of compiler flag	Version of GCC compiler				
		4.6.4	4.7.4	4.8.5	4.9.3	5.1.0
-O0	1					
	2					
	3	(1)				
	4					
	5					
	6					
	7					
	8					
-O1	1					
	2					
	3					
	4					
	5					(2)
	6					
	7					
	8					
-O2	1					
	2					
	3					
	4					
	5					(3)
	6					
	7					
	8					
-O3	1					
	2					
	3	(4)	(5)	(6)	(7)	(8)
	4					
	5					(9)
	6					
	7					
	8					

2

1 Table 13: Similar to Table 11 except for the simulation results of FGOALS-g2. FGOALS-g2
 2 has not been compiled using the GCC compilers for simulation runs before. Therefore a large
 3 proportion of simulation runs are failed (marked with "--" in the table). For example, crashes or
 4 deadlocks are encountered under compiler optimization levels O1 to O3.

Compiler optimization level	No. of compiler flag	Version of GCC compiler				
		4.6.4	4.7.4	4.8.5	4.9.3	5.1.0
-O0	1					
	2					
	3	(1)		(2)		
	4					
	5	(3)	(4)			
	6					
	7					
	8					
-O1	1	--	--	--	--	--
	2	--	--	--	--	--
	3	--	--	--	--	--
	4	--	--	--	--	--
	5	--	--	--	--	--
	6	--	--	--	--	--
	7	--	--	--	--	--
	8	--	--	--	--	--
-O2	1	--	--	--	--	--
	2	--	--	--	--	--
	3	--	--	--	--	--
	4	--	--	--	--	--
	5	--	--	--	--	--
	6	--	--	--	--	--
	7	--	--	--	--	--
	8	--	--	--	--	--
-O3	1	--	--	--	--	--
	2	--	--	--	--	--
	3	--	--	--	--	--
	4	--	--	--	--	--
	5	--	--	--	--	--
	6	--	--	--	--	--
	7	--	--	--	--	--
	8	--	--	--	--	--

5

1 Table 14: Ideal bitwise identical compiling setup sets of the three models when using Intel
 2 compilers. Each color except the white corresponds to an ideal bitwise compiling setup set.

Compiler optimization level	No. of compiler flag	Version of Intel compiler				
		11	12	13	14	15
-O0	1	Yellow	Red	Red	Red	Red
	2	Yellow	Red	Red	Red	Red
	3	White	Red	Red	Red	Red
	4	Yellow	Red	Red	Red	Red
	5	Yellow	Red	Red	Red	Red
	6	Yellow	Red	Red	Red	Red
	7	Yellow	Red	Red	Red	Red
	8	Yellow	Red	Red	Red	Red
	9	Yellow	Red	Red	Red	Red
-O1	1	Yellow	Red	Red	Red	Red
	2	Yellow	Red	Red	Red	Red
	3	White	White	White	White	White
	4	Yellow	Red	Red	Red	Red
	5	Yellow	Red	Red	Red	Red
	6	Yellow	Red	Red	Red	Red
	7	Yellow	Red	Red	Red	Red
	8	Yellow	Red	Red	Red	Red
	9	Yellow	Red	Red	Red	Red
-O2	1	Yellow	Red	Red	Red	Red
	2	Yellow	Red	Red	Red	Red
	3	White	White	White	White	White
	4	Yellow	Red	Red	Red	Red
	5	Yellow	Red	Red	Red	Red
	6	Yellow	Red	Red	Red	Red
	7	Yellow	Red	Red	Red	Red
	8	Yellow	Red	Red	Red	Red
	9	Yellow	Red	Red	Red	Red
-O3	1	Yellow	Red	Red	Red	Red
	2	Yellow	Red	Red	Red	Red
	3	White	White	White	White	White
	4	Yellow	Red	Red	Red	Red
	5	Yellow	Red	Red	Red	Red
	6	Yellow	Red	Red	Red	Red
	7	Yellow	Red	Red	Red	Red
	8	Yellow	Red	Red	Red	Red
	9	Yellow	Red	Red	Red	Red

3

4

1 Table 15: Examples of different results of the calculation at line 330 of Fig. 2 when changing
 2 the compiler optimization level from *O1* to *O2*. The input of the calculation is the same (bitwise
 3 identical) at both compiler optimization levels. The different digits in the results are highlighted
 4 in red.

Variables		Example			
		No. 1	No. 2	No. 3	
I n p u t	k	2	3	3	
	i	9	1	5	
	ipair	1	1	1	
	dryvol_t_new(ipair,i,k)	1.245177471001780E-013	1.367964902074264E-013	1.362619492656580E-013	
	num_t_oldbnd(ipair,i,k)	660367763.850537	673856916.583178	665467981.351062	
	factoraa(mfrm)	1.41486733199200	1.41486733199200	1.41486733199200	
O u t p u t	dgn_t_new(ipair,i,k)	optimization level O1	5.107909846347498E-008	5.235166698430766E-008	5.250216806732902E-008
		optimization level O2	5.107909846347492E-008	5.235166698430762E-008	5.250216806732898E-008

5

1 Table 16: Assembly codes of the calculation at line 330 in Fig. 2 in two compiler optimization
 2 levels (*O1* and *O2*). The most significant difference of the assembly codes is the calling of
 3 different power functions.

Optimization level O1	Optimization level O2
movq %rsi, -40(%rbp)	movsd %xmm7, -344(%rbp)
movq %r8, -32(%rbp)	movsd %xmm1, -320(%rbp)
movq %r9, -24(%rbp)	movsd %xmm3, -312(%rbp)
movsd %xmm8, -16(%rbp)	movsd %xmm2, -304(%rbp)
call pow	call cbrt

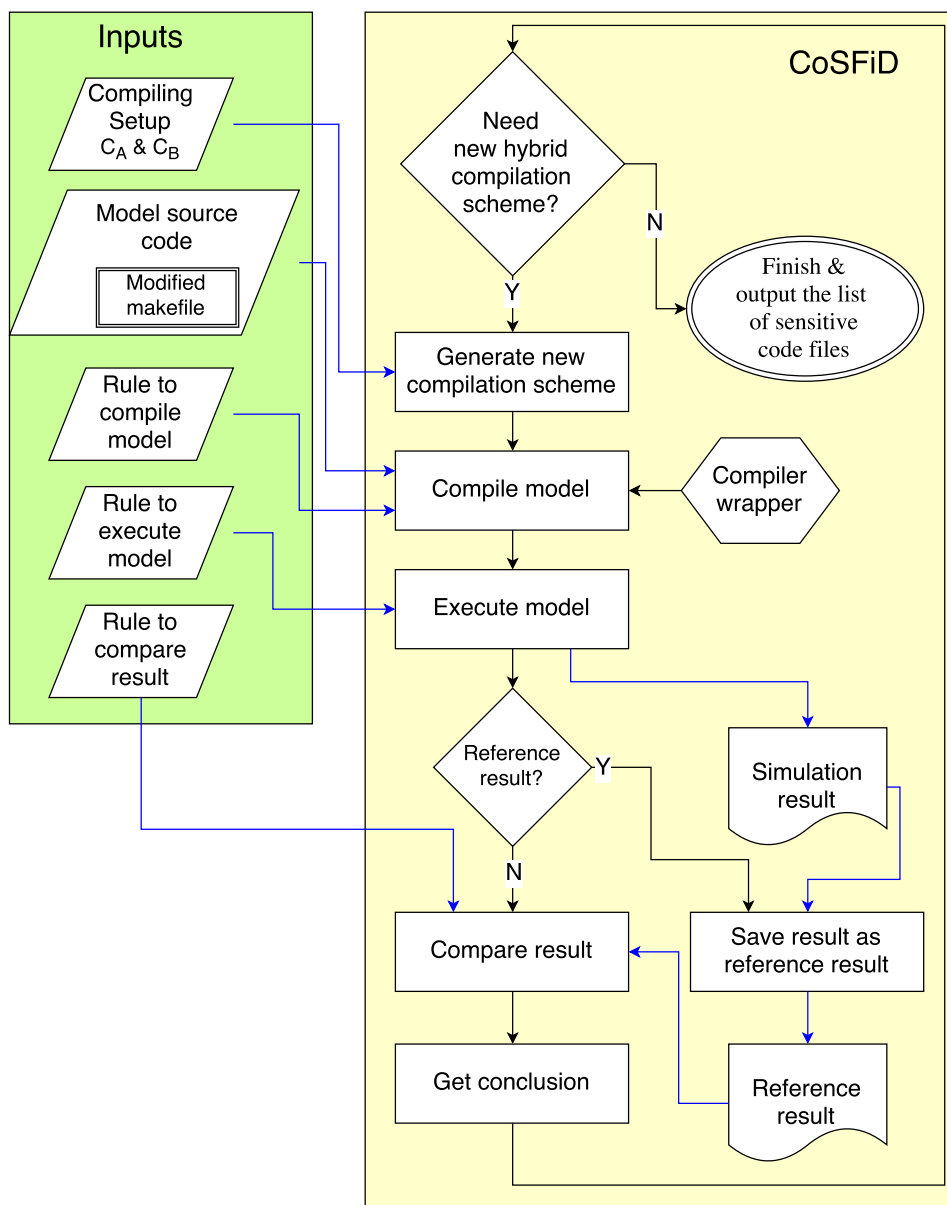
4

5

1 Table 17: An example of obvious different results in lines 1923-1932 of Fig. 3 when changing
 2 the compiler optimization levels (from O2 to O3). A manual result calculated by Python is also
 3 provided.

Variables		Value	
Input	WORK3(i, j)	0.0000000000000000	
	dz(k)	1000.000000000000	
	KAPPA_ISOP(i, j, kbt, k, bid)	1.991793396882581E-006	
	SLX(i, j, ieast, kbt, k, bid)	-0.120114394362605	
	HYX(i, j, bid)	1.32933181234324	
	TX(i, j, k, n, bid)	-0.161989629268646	
	SLY(i, j, jnorth, kbt, k, bid)	-0.13980933 0009777	
	HXY(i, j, bid)	0.761427260631393	
	TY(i, j, k, n, bid)	0.152039408683777	
	SLX(i, j, iwest, kbt, k, bid)	-7.344871974748127E-002	
	HYX(i-1, j, bid)	1.32933181234324	
	TX(i-1, j, k, n, bid)	-0.192202091217041	
	SLY(i, j, jsouth, kbt, k, bid)	-0.112685940441552	
	TY(i, j-1, k, n, bid)	0.743088001419362	
	TY(i, j-1, k, n, bid)	0.122525990009308	
Output	WORK3(i, j)	Execution result (at optimization level O2)	3.622331248054413E-005
		Execution result (at optimization level O3)	3.571897176404182E-005
		Manual result (by Python)	3.62233124805436E-05

4



1
 2 Figure 1: Flowchart of CoSFID for detecting compilation-sensitive code files. In each iteration,
 3 CoSFID first checks whether it is necessary to generate a new hybrid compilation scheme (some
 4 code files are compiled with C_A and the remaining code files are compiled with C_B). If
 5 unnecessary, which means the whole process of the detection should end, CoSFID will output
 6 all compilation-sensitive code files. Otherwise, CoSFID generates a new hybrid compilation
 7 scheme, and then calls the corresponding rule to compile the model code using the compiler
 8 wrapper and run the simulation. If it is the first run of the simulation, which also means all code
 9 files are compiled with C_A , the simulation result will be recorded as the reference result.
 10 Otherwise, CoSFID calls the corresponding rule to compare the simulation result to the
 11 reference result and then uses the conclusion to drive the next iteration.

```

321 ! num_t_old is total number in particles/kmol-air
322   num_t_old = q(i,k,numptr_amode(mfrm)-loffset)
323   num_t_old = num_t_old + qqcw(i,k,numptrcw_amode(mfrm)-loffset)
324   num_t_old = max( 0.0_r8, num_t_old )
325   dryvol_t_oldbnd = max( dryvol_t_old, dryvol_smallest(mfrm) )
326   num_t_oldbnd = min( dryvol_t_oldbnd*v2nlorlx(mfrm), num_t_old )
327   num_t_oldbnd = max( dryvol_t_oldbnd*v2nhirlx(mfrm), num_t_oldbnd )
328
329 ! no renaming if dgnum < "base" dgnum,
330   dgn_t_new = (dryvol_t_new/(num_t_oldbnd*factoraa(mfrm)))*onethird
331   if (dgn_t_new .le. dgnum_amode(mfrm)) cycle mainloop1_ipair
332
333 ! compute new fraction of number and mass in the tail (dp > dp_cut)
334   lndgn_new = log( dgn_t_new )
335   lndgv_new = lndgn_new + dum3alnsg2(ipair)
336   yn_tail = (lndp_cut(ipair) - lndgn_new)*factoryy(mfrm)
337   yv_tail = (lndp_cut(ipair) - lndgv_new)*factoryy(mfrm)
338   tailfr_numnew = 0.5_r8*erfc( yn_tail )
339   tailfr_volnew = 0.5_r8*erfc( yv_tail )

```

1

2

3 Figure 2: Part of the code lines of the compilation-sensitive code segment in the code file
4 *modal_aero_rename.F90* of CAM5. It is found that the code at line 330 can produce different
5 results when different compiling setups are used.

6


```

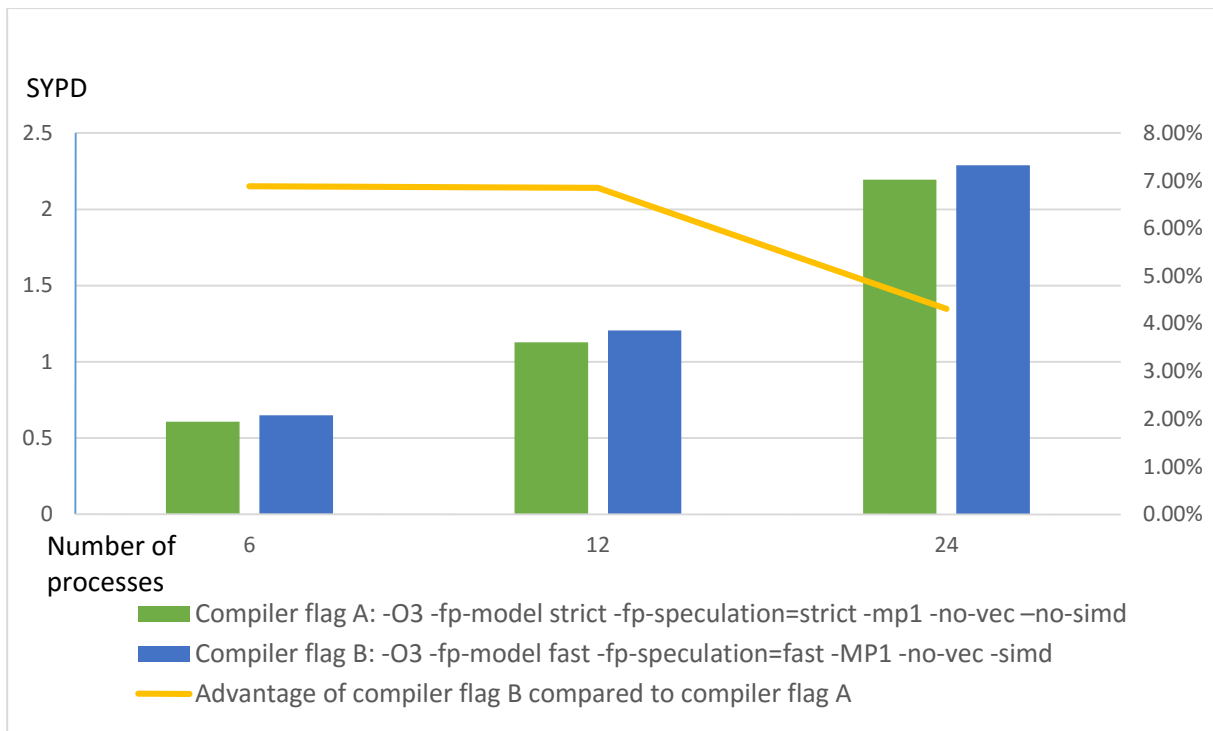
1920      do j=this_block%jb,this_block%je
1921          do i=this_block%ib,this_block%ie
1922
1923              WORK3(i,j) = WORK3(i,j) &
1924                  + ( dz(k) * KAPPA_ISOP(i,j,kbt,k,bid) &
1925                      * ( SLX(i,j,ieast,kbt,k,bid) &
1926                          * HYX(i,j,bid) * TX(i,j,k,n,bid) &
1927                          + SLY(i,j,jnorth,kbt,k,bid) &
1928                          * HXY(i,j,bid) * TY(i,j,k,n,bid) &
1929                          + SLX(i,j,iwest,kbt,k,bid) &
1930                          * HYX(i-1,j,bid) * TX(i-1,j,k,n,bid) &
1931                          + SLY(i,j,jsouth,kbt,k,bid) &
1932                          * HXY(i,j-1,bid) * TY(i,j-1,k,n,bid) ) ) &
1933
1934          enddo
1935      enddo
1936
1937      do j=this_block%jb,this_block%je
1938          do i=this_block%ib,this_block%ie
1939
1940              WORK3(i,j) = WORK3(i,j) &
1941                  + ( SF_SLX(i,j,ieast,kbt,k,bid) &
1942                      * HYX(i,j,bid) * TX(i,j,k,n,bid) &
1943                      + SF_SLY(i,j,jnorth,kbt,k,bid) &
1944                      * HXY(i,j,bid) * TY(i,j,k,n,bid) &
1945                      + SF_SLX(i,j,iwest,kbt,k,bid) &
1946                      * HYX(i-1,j,bid) * TX(i-1,j,k,n,bid) &
1947                      + SF_SLY(i,j,jsouth,kbt,k,bid) &
1948                      * HXY(i,j-1,bid) * TY(i,j-1,k,n,bid) ) &
1949
1950          enddo
1951      enddo

```

1

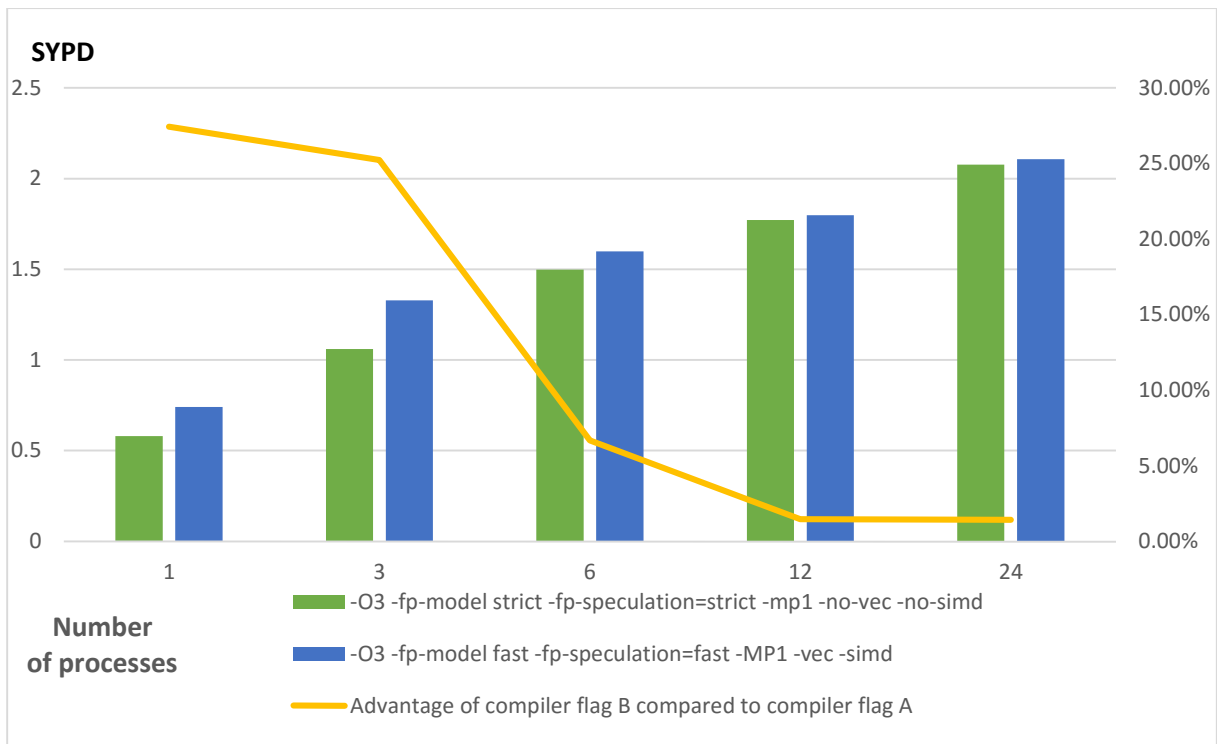
2 Figure 3: Part of the code lines of the compilation-sensitive code segment in the code file
3 *hmix_gm.F90* of POP2. It is found that the code from line 1923 to line 1932 can produce
4 significantly different results when different compiling setups are used.

5



1
2
3
4
5
6
7
8
9
10
11
12
13
14

Figure 4: Simulation speed (simulated years per day; SYPD) of CAM5 under two compiler flags (A and B) of Intel compiler version 13 when increasing the number of processes from 6 to 24. The high-performance computer Tansuo100 is used for this test. Compiler flag A (“-O3 -fp-model strict -fp-speculation=strict -mp1 -no-vec -no-simd”) is from the biggest bitwise identical compiling setup sets in Table 8. Compiler flag B (“-O3 -fp-model fast -fp-speculation=fast -MP1 -no-vec -simd”) should be the compiler flag for fastest simulation speed. Compiler flag “-O3 -fp-model fast -fp-speculation=fast -MP1 -vec -simd” should be more aggressive than compiler flag B in compiler optimizations. It is not used in this test because the corresponding simulation run of CAM5 crashes. The advantage of compiler flag B compared to compiler flag A is defined as the performance improvement when compiler flag is changed from A to B.



1

2 Figure 5: Simulation speed (simulated years per day; SYPD) of GAMIL2 (Li et al, 2013b) under
 3 two compiler flags (A and B) of Intel compiler version 13 when increasing the number of
 4 processes from 1 to 24. The high-performance computer Tansuo100 is used for this test.

5 Compiler flag A (“-O3 -fp-model strict -fp-speculation=strict -mp1 -no-vec -no-simd”) is the
 6 also the compiler flag A used in Fig. 4. Compiler flag B (“-O3 -fp-model fast -fp-
 7 speculation=fast -MP1 -vec -simd”) should be the compiler flag for fastest simulation speed.

8 The advantage of compiler flag B compared to compiler flag A is defined as the performance
 9 improvement when compiler flag is changed from A to B.