# MP-CodeCheck User's Manual (v.0.0.6)

Merly, Inc. (Date: 2022-05-28)

## ABSTRACT

Software debugging, the process of identifying and mitigating software defects, has been reported to consume upwards of 50% of all software development time. Software defects can arise from a number of issues. These can include logical errors, temporally or spatially inefficient code, or poorly designed and implemented code, also known as technical debt, to name a few. Some of the most severe software defects are security vulnerabilities, which can compromise the safety of a software system, an organization, or even an entire company. We created MP-CodeCheck (MPCC) to help programmers identify and address such defects (also known as anomalies). In this guide, we walk users through MPCC's general uses as well as the steps required to setup, run, and perform inference against a code base.

*Your code is never sent to us. Ever. Your privacy is our #1 priority.*

## Contents

# Merly's MP-CodeCheck

## 1 What is MP-CodeCheck?

MP-CodeCheck (MPCC) is an AI-based code anomaly detection system. More specifically, MPCC uses self-supervision, federated learning, and programmatic-guided evolution to detect anomalous code patterns. An overview of its system design is shown in Figure 1. MPCC was designed to learn good and bad code syntax, patterns, and semantics from a large corpora of existing code. Once trained, MPCC's model can be used for a variety of tasks such as: *(i)* detecting potential anomalies in existing code, *(ii)* grading the quality of an existing repository, and *(iii)* guiding programmers through the important aspects of a new code repository, to name a few.
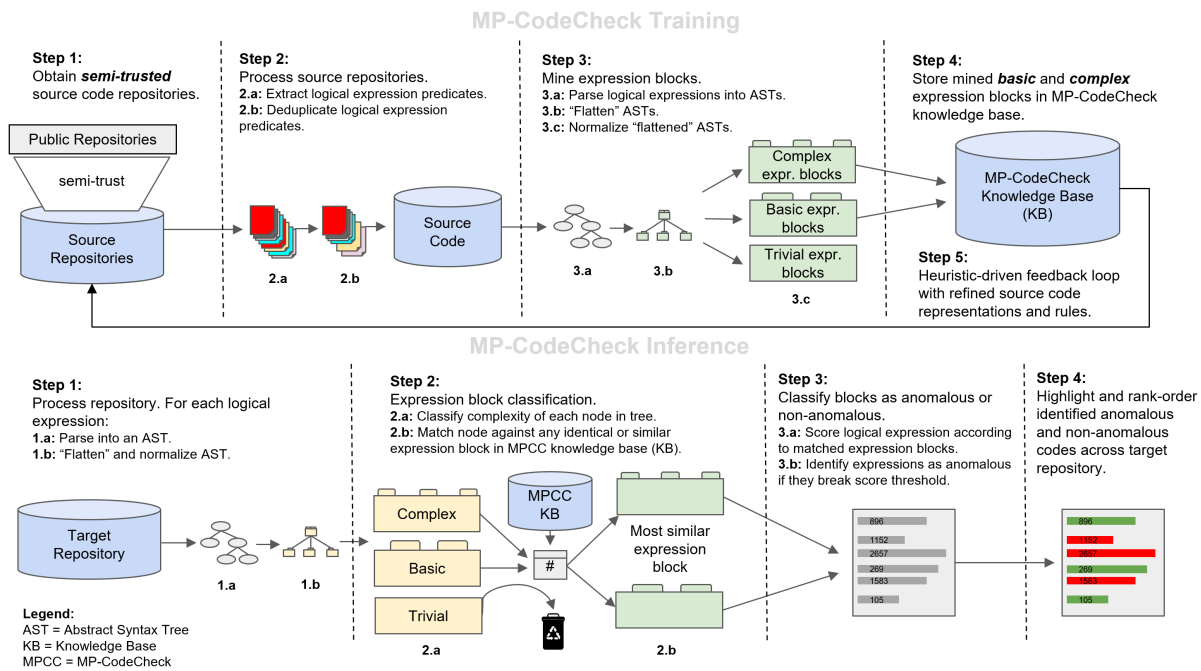


**Figure 1.** System Overview of MP-CodeCheck.

At its core, MPCC is a machine programming (MP) system that autonomously identifies anomalous logical expressions directly in source code. These anomalous expressions, also known as anomalies, are often latent defects in the existing code that programmers have failed to identify or correct. MPCC helps programmers find these anomalies and correct them, thereby improving the overall quality of the existing software. For this limited release version of MPCC, we only include MPCC the ability to perform inference (i.e., detect good or bad patterns) on code. In subsequent releases of MPCC, we may also include the ability to train new models on other code bases, including users' own proprietary ones.

## 2 Installation Instructions

Below we list MPCC's installation instructions for the currently supported operating systems (OSes). If you previously installed MPCC without a product key, but now have one, it is safe to run the MerlyInstaller again to register MPCC with a product key. A product key is a 16 character string, separated by hyphens each four characters and is generated when you register MPCC on Merly's website (e.g., 5SA9-HBP2-WRBV-5WA1).

If you encounter any trouble with these steps, please contact `support@merly.ai` for assistance.

### 2.1 Linux (CentOS, RedHat, SUSE, Ubuntu), MacOS (M1 ARM, x64 Intel)

Launch a command line interface (CLI) and execute the following command:

```
/bin/bash -c "$(curl -fsSL https://raw.githubusercontent.com/merly-ai/MPCC-Universal/main/install.sh)"
```

If you have not registered with a product key, execute the following command from the CLI, where `<key>` is your product key:

```
./MerlyInstaller -k <key> install
```

For MacOS, we recommend running MPCC with iTerm2, due to its support of a broader color scheme than is possible for the default MacOS terminal. You can download it for free here: https://iterm2.com/downloads.html.

## 2.2 Windows (64-bit)

Launch `cmd.exe` (do not use PowerShell as the below cURL command will not work). Navigate to your user preferred installation directory (e.g., `cd C:\Users\Paul`). Then execute the following commands where `<key>` is your product key:

```
mkdir MPCC
cd MPCC
curl -LO https://github.com/merly-ai/MP-CodeCheckBin-Windows/raw/main/bin/latest/MerlyInstaller.exe
MerlyInstaller -k <key> install
```

## 2.3 Updating MPCC with the MerlyInstaller

MPCC is constantly being updated. To update your local copy of MPCC to the latest version including updating all of latest programming language models it supports, simply run the MerlyInstaller (installed in the above installation steps) in the following way:

```
MerlyInstaller updateall
```

If you would like to update only certain components, please use the 'usage' command from the MerlyInstaller as shown below for more details on how to update only the components you are interested in.

```
MerlyInstaller usage
```

# 3 Launching MP-CodeCheck

Now that setup is complete, let's launch MPCC to perform inference analysis. From the command line interface (CLI), type the following (where "[code base folder]" is a directory that contains the code you want to analyze):

**MacOS, Linux** `./MPCC infer -D [code base folder]`

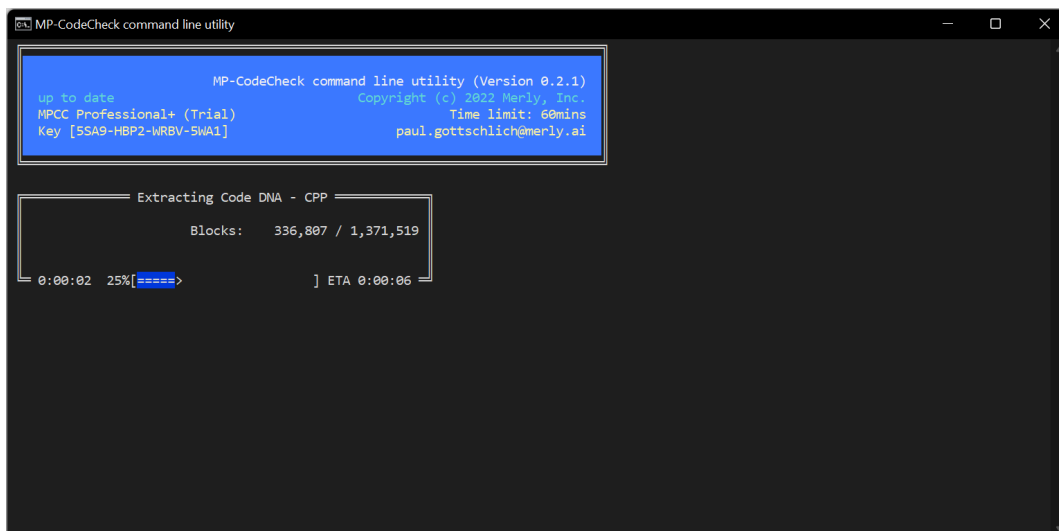**Windows** `MPCC.exe infer -D [code base folder]`



**Figure 2.** Launching MP-CodeCheck and Extracting Code DNA.

When run successfully, MPCC will display information that looks similar to the screenshot shown in Figure 2. This shows the progress of MPCC extracting the code DNA from the training model.
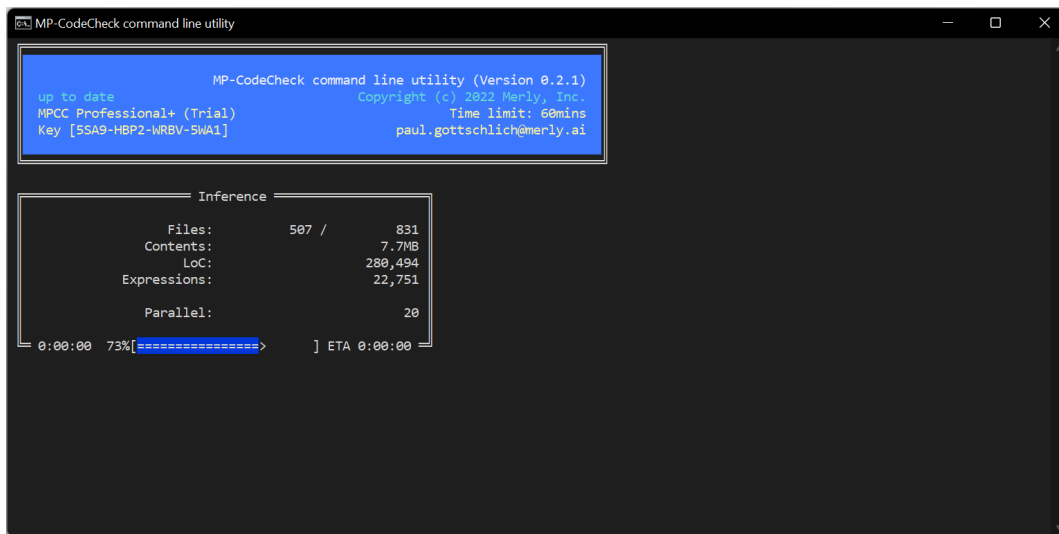
**Figure 3.** MP-CodeCheck Performing Inference Processing.

When MPCC has loaded its trained model and processed the code DNA, it begins inference analysis on all source code that it finds in the files of the directory (or subdirectories) you have supplied when launching it. Figure 3 shows an example of MPCC's inference progress in analyzing a code repository, how much work it has completed, and how much work is remaining. When inference analysis has completed, the *Code View* screen will appear (shown in Figure 4), which will allow a user to analyze the inference results as discussed in the next section.

## 4 Exploring MPCC's Inference Results

After inference analysis is performed, MPCC will show a user interface that includes source code, with an expression highlighted. We call this screen the *Code View*, which will be described in more detail in Views section of this manual. Figure 4 provides an example of an anomalous code example found by MPCC.

**Figure 4.** An Anomalous Example in MP-CodeCheck's Code View.



**Figure 5.** A Non-Anomalous Example in MP-CodeCheck's Code View.

**Sort Criteria:**   This refers to how MPCC is sorting the list of expressions it has found. This can be via score (a numeric value assigned by anomaly identification and complexity), or location (sequential code order).

**Class Filter:**   This refers to which class of complexity is being filtered in the current view. This can be set from a minimum value of trivial to a high value of Max complexity.

**Cost Filter:**   This refers to a "mental cost" of an expression. This filter can be set from a minimum value of 0 to a maximum value of 2,000.

**Displayed Items:**   This refers to which items MPCC is displaying. It can be set to all expressions, or only anomalous expressions.

**Hide/Show Known Good:**   This refers to whether or not MPCC displays expressions that have been marked by the user as Known Good.

**Anomaly Identification:**   This displays whether or not MPCC has identified the current expression as an anomaly. Non-anomalous expressions will be classified as "known pattern detected" and highlighted in green. Anomalous expressions will be classified as "unfamiliar pattern(s) detected" and will be highlighted in green.

**Cost:**   This displays the "mental cost" of the current expression.

**Complexity:**   This displays the class of complexity of the current expression.

**Source Code Location:**   This displays the file location of the source code under review.

**Anomaly/Expression Count:**   This displays the count of the highlighted expression, as well as the total expressions found in this file. If the user toggles the filter to show only anomalies, this will display the count of highlighted anomaly, and the total anomalies found in the current file.

**Walking Through Code:**   You can move forwards and backwards through the expressions by using the left and right arrow keys, and can page up and page down through the code (by location) using the Page Up and Page Down keys. You can also scroll up and down through the code by hold the Control key while pressing the up or the down arrow, respectively.

## 5  Basic Commands and Views

In MPCC, there are a number of supported keyboard and mouse commands. In this section we describe those keystrokes and explain mouse behavior. Perhaps the most important initial command to remember is the *help* command which can be launched by pressing the character 'h' on your keyboard. The help command lists all of the keyboard commands, so if you ever find yourself not remembering a keyboard command, just press 'h' and MPCC will launch the keyboard shortcut commands. A screenshot of the help dialogue box is shown in Figure 11.

In addition to commands, there are several screens users can utilize to help them gain deeper insights into specific anomalies, general anomaly information, anomalies by file, anomalies per file, and so forth.

**Figure 6.** MP-CodeCheck's Code View.

**Code View:** This is the view of all of the code, with the expressions found highlighted. This view is the default view when MPCC is initially run.



**Figure 7.** MP-CodeCheck's Anomalies View.

**Anomalies View:** Press 'a' to switch to the Anomalies view. This view shows all of the expressions in the code (across all files) that MPCC has determined to be an anomaly, sorted by score. You can move up and down the list using the up and down arrows, or the Page Up and Page Down keys. Press Enter with an anomaly highlighted to switch back to the Code View of that specific anomaly.
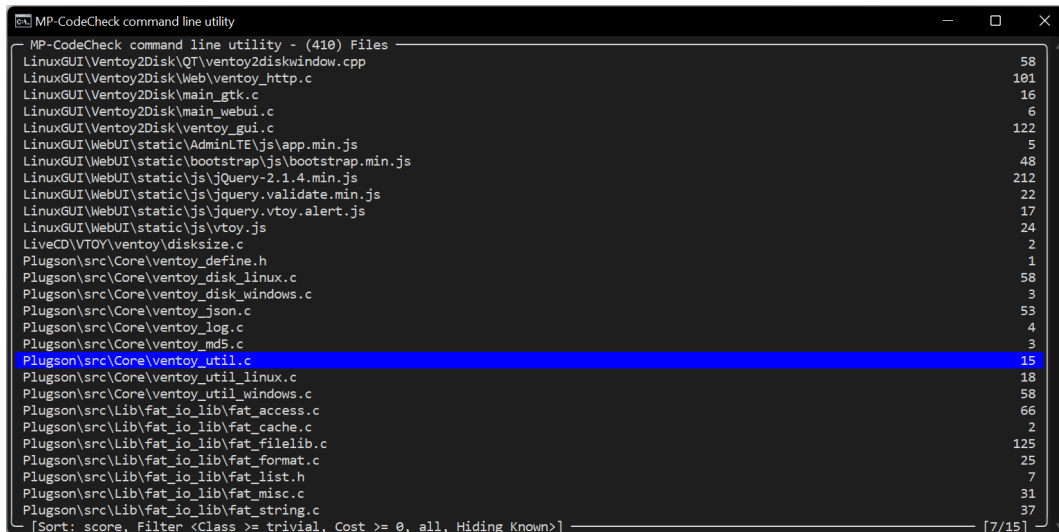
**Figure 8.** MP-CodeCheck's Files View.

**Files View:** Press 'f' to switch to the Files view. This view shows all of the source code files, with the total number of expressions MPCC found in each file. You can move up and down the list using the up and down arrows, or the Page Up and Page Down keys. Press Enter with a file highlighted to switch back to the Code View of the expressions within that specific file.
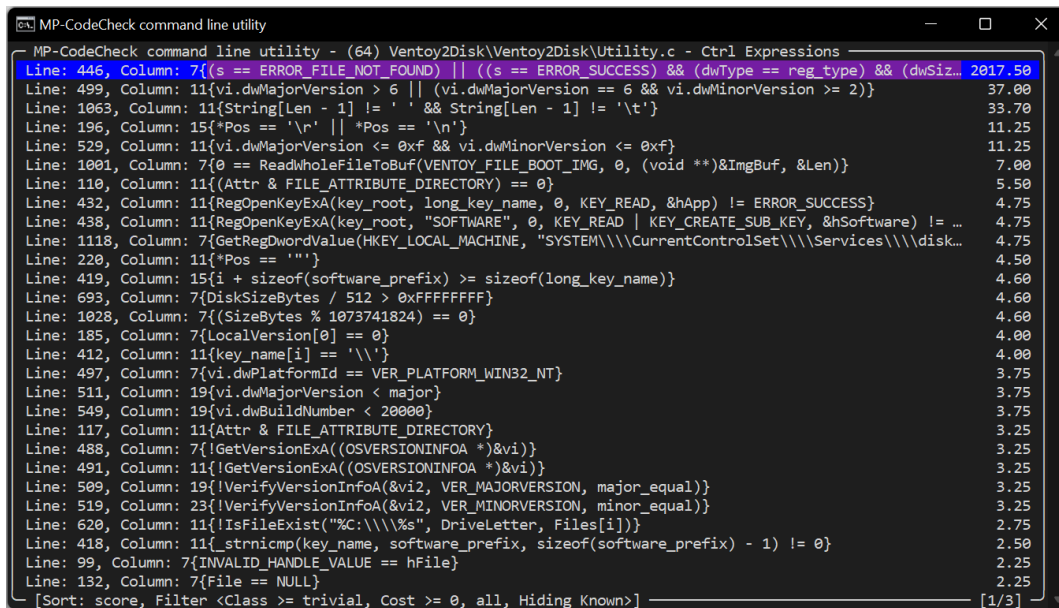


**Figure 9.** MP-CodeCheck's Expressions View.

**Expressions View:** Press 'e' to switch to the Expressions view. This view shows all of the expressions in the current file, sorted by score. You can move up and down the list using the up and down arrows, or the Page Up and Page Down keys. Also note that you can toggle the sort between code location and score by pressing the 's' key. Press Enter with an expression highlighted to switch back to the Code View with that specific expression highlighted.

**Figure 10.** MP-CodeCheck's Details View.

**Details View:** Press 'd' to switch to the Details view. This view shows the detail of the currently selected expression. The detail lets you know how many anomalies MPCC identified within the expression, the cost, and the total score. Press 'd' to return to Code View.

**Figure 11.** MP-CodeCheck's Help Dialogue Box.

**Help Pop-up:**   In addition to the above views, you can press the 'h' key in any view to bring up the help screen which will show you all of the hot keys and their functions.

# 6  Sorting/Filtering Inference Results

The following lists the ways MPCC's inference results on source code data can be sorted and/or filtered.

*Sort Criteria:*

**Options:**

- Score (numeric value assigned by anomaly identification and complexity)

- Location (sequential code order)

**Default:** Score
**Toggle:** 's' key

*Class filter:*

**Options:**

- Trivial (minimum)

- Basic

- Complex 1

- Complex 2

- Max

**Default:** Trivial
**Adjust:** '1', '2', '3', '4', '5' keys

### *Cost filter:*

**Options:**

- 0 (minimum) to 2,000

**Default:** 0
**Adjust:** ',' to decrease, '.' to increase, 'm' to reset to 0 (minimum)

### *Displayed items:*

**Options:**

- All expressions

- Anomalies Only

**Default:** All expressions
**Toggle:** '0' key

### *Hide/Show Known Good:*

**Options:**

- Hide Known Good

- Show Known Good

**Default:** Hide Known Good
**Toggle:** '9' key

## MPCC Generated Files

In addition to the live (online) user interface, you can also review the inference results offline through four MPCC generated files. These files are re-generated each time inference is run successfully. These files will be created in the same folder that the MPCC executable was launched and have the following naming structure.

**[Code Repo].by_file.txt:**   This file lists all anomalous expressions (that are not nested if's) found by MPCC. This human readable file lists the original anomalous source and its normalized version.

**[Code Repo].by_file_nested_if.txt:**   This file lists all nested if expressions that are found by MPCC to be anomalous. This human readable file lists the original anomalous source and its normalized version.

**[Code Repo].mpcc.anomaly_list.json:**   This file lists all expressions that are found by MPCC to be anomalous, in a machine-readable format.

**[Code Repo].mpcc.summary.json:**   This file contains a summary of all of the files, size, and lines of code reviewed by MPCC. It also provides a summarized report of the number of expressions, anomalies, and scores found in the source code that inference was performed on, in a machine-readable format.

# 7 Advanced Usage

In this section, we describe how to use some of the advanced features of MP-CodeCheck.

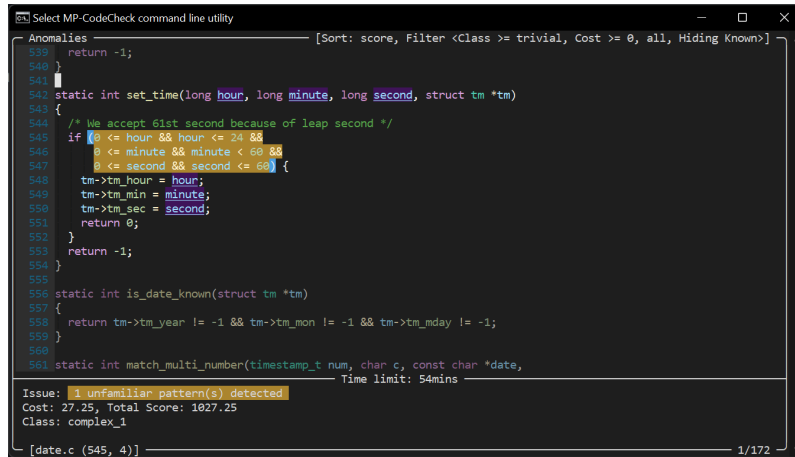## 7.1 Copy/Paste Functionality in MPCC



**Figure 12.** A white cursor appears to mark the start of your copy area.

Left-click on the beginning of the area you'd like to copy. You'll see a white cursor where the copy will begin (Figure 12). Next, hold down the left mouse button and drag the cursor over the area you'd like to copy, as shown in Figure 13). You can also hold down the Shift key and use the cursor keys (up/down, left/right) to select the area you'd like to copy. When you have the selection highlighted, you can either right-click on your mouse or press the Enter key to complete the copy.
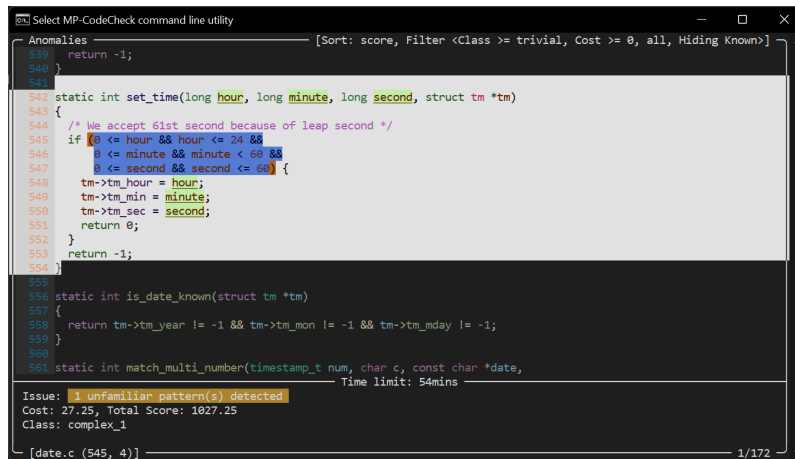


**Figure 13.** Highlighted text is ready to be copied.

## 7.2 MPCC Configuration

For users who wish to customize their MPCC experience, a JSON file is available to configure MPCC to fit your preferences.

The JSON file is located at the following location:

```
%appdata%\..\local\merly.ai\debugging\MP-CodeCheck\config.json
```

You can use any text editor to modify the colors, log file locations, and settings. Let's take a closer look.

Colors: These are stored in the JSON file in hexadecimal (HEX) RGB; simply use your favorite color picker to find the hex value of the color you'd like, and change the value of the associated item.

For example, you can set `anomaly_background` to RGB `ab852e` to change the highlight color of the anomalous expressions to dark orange.
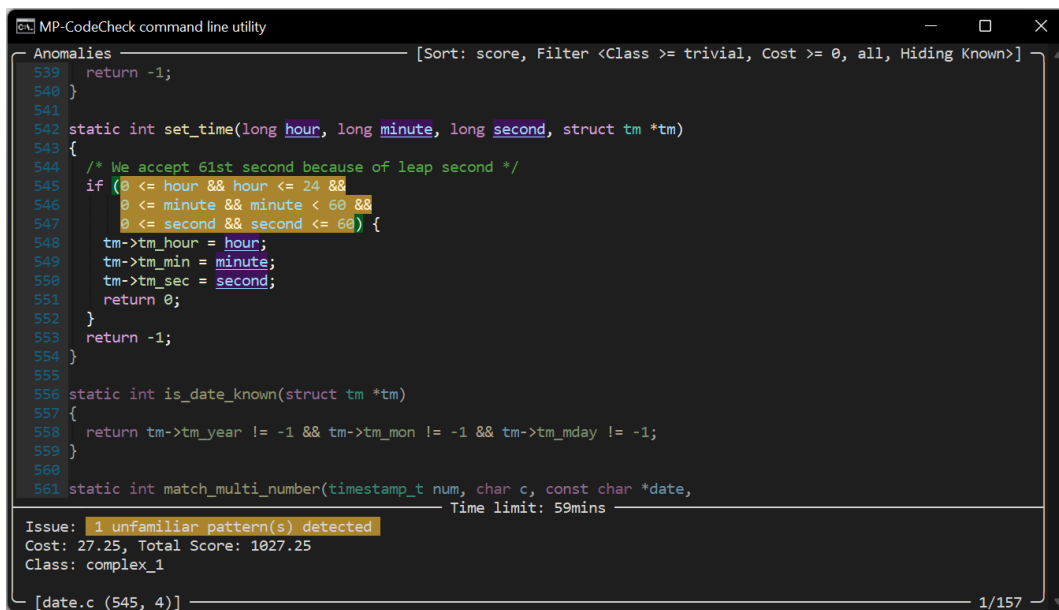


**Figure 14.** An example of changing the anomaly highlight color in MP-CodeCheck's Code View.

Or, set `highlight_background` to RGB `4a9de0` to change the highlight color of the non-anomalous expressions to light blue.

Models Path: The models path can be modified in the JSON file, and is initially set to:

```
"models_path": "[initial install directory]\\models"
```

Log Path: The log path can be modified in the JSON file, and is initially set to:

```
"log_path": "[initial install directory]\\logs"
```

Settings:
"run_training" – Determines whether or not training should be run before inference on the source code (defaults to true).
"filter" – Determines if items (such as nested_ifs) are extracted (defaults to true).

**Figure 15.** An example of changing the expression highlight color in MP-CodeCheck's Code View.

### 7.3 File and Folder Exclusion

If you wish to omit files or folders from MP-CodeCheck's inference, you can use MPCC's file and folder exclusion command. This feature is run during the inference step and is triggered by the -X command for directories (uppercase X) and -x for files (lowercase X). As an example, one might trigger it using the following format:

**MacOS, Linux**  `./MPCC infer -D [code directory] -X [directories to exclude] -x [files to exclude]`

**Windows**  `MPCC.exe infer -D [code directory] -X [directories to exclude] -x [files to exclude]`

MPCC's file and directory exclusion command uses wildcards. The following are some example usages:

- `-x *api*` (lowercase x, excludes all files with 'api' in the name)

- `-X */out/*` (uppercase X, excludes all folders that contain a folder 'out' in the name)

- `-x */internal` (lowercase x, excludes all files with 'internal' as the last folder name)

For more details launch MPCC with the 'usage' command.

THE MACHINE PROGRAMMING COMPANY