

MP-CodeCheck User's Manual (v.0.0.3)

Merly, Inc.

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. We discuss MPCC's core features around code anomaly inference and how to leverage them in various views inside of MPCC and outside of it.



Merly's MP-CodeCheck

What is MP-CodeCheck?

MPCC is an AI-based code anomaly detection system. More specifically, MPCC uses self-supervision, iterative 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.

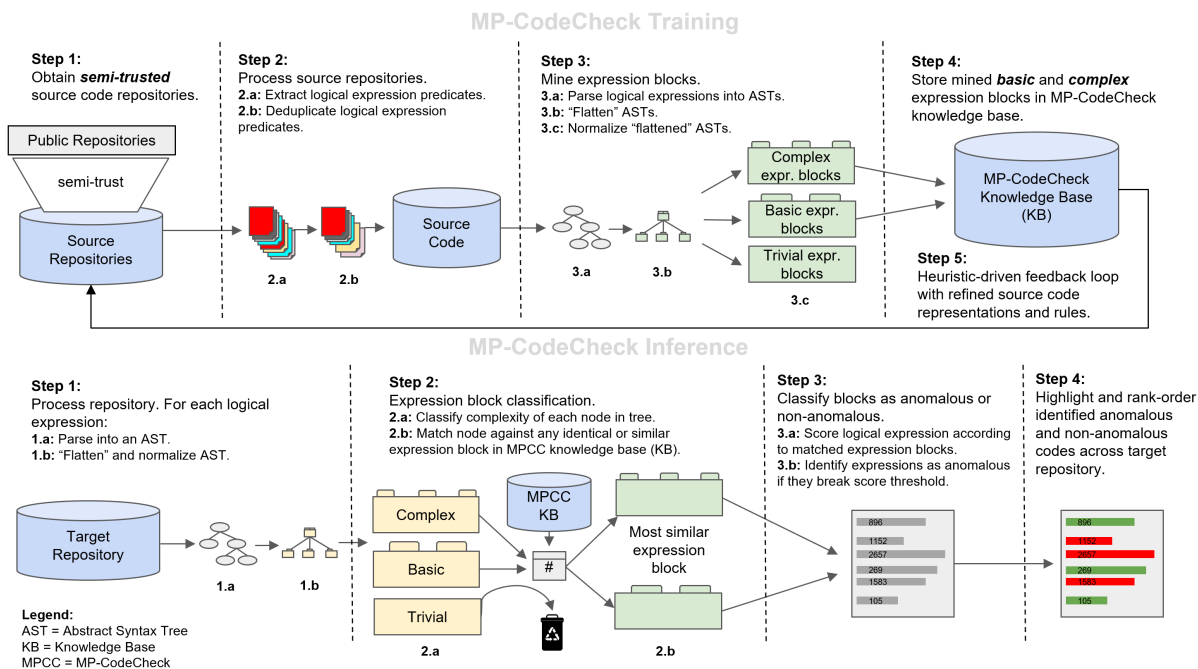


Figure 1. System Overview of MP-CodeCheck.

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.

Pre-Setup Instructions

Merly will provide you with two files that are necessary to run MPCC: (i) its model file and (ii) the MPCC executable. These files are large, so we zip them using 7Zip. You can use this tool (or your own preferred 7zip compression tool) to unzip the contents.

You can download the 7Zip tool here: <https://www.7-zip.org/>

Setup Instructions

Prior to running inference and reviewing the results, let's set up the environment. To run MPCC, you'll need the following three things (at a minimum):

1. A model trained on code (provided in the Merly zip file).
2. The MPCC executable (provided in the Merly zip file).
3. A code base to run inference against (provided by you, the user).

Once you have unzipped the file provided by Merly, you will notice two main files:

- The model: `model.database.bin`
- The executable: `MPCC.exe`

Please ensure you place both the MPCC model and the executable files in the same folder. Then, to simplify inference, we recommend you place the code repository folder in the same directory as MPCC. Your setup is now complete!

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):

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



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

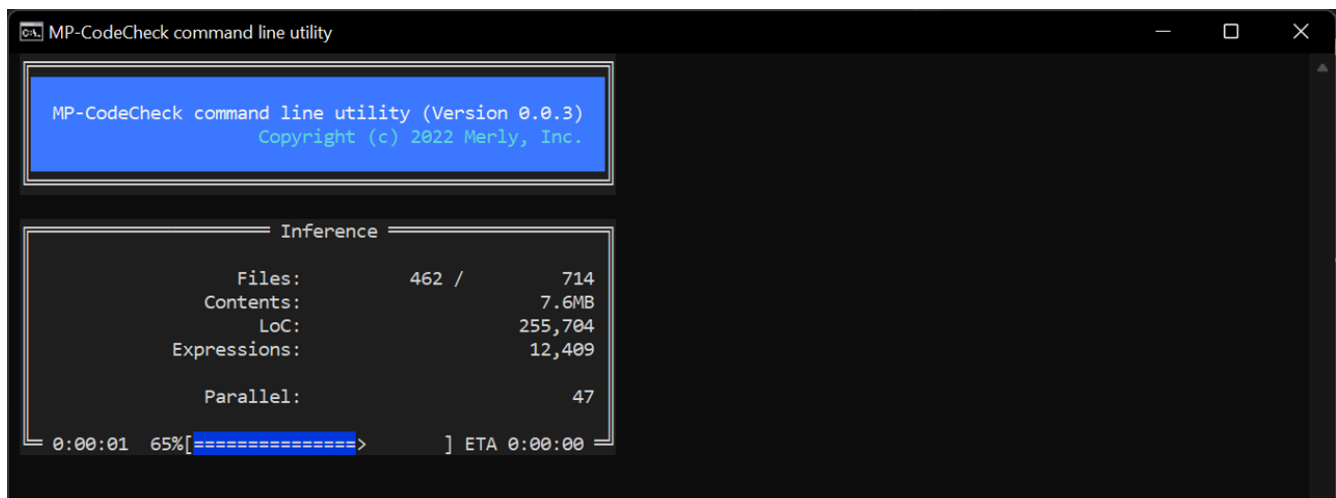


Figure 3. MP-CodeCheck Performing Inference Processing.

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.

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.

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.

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.

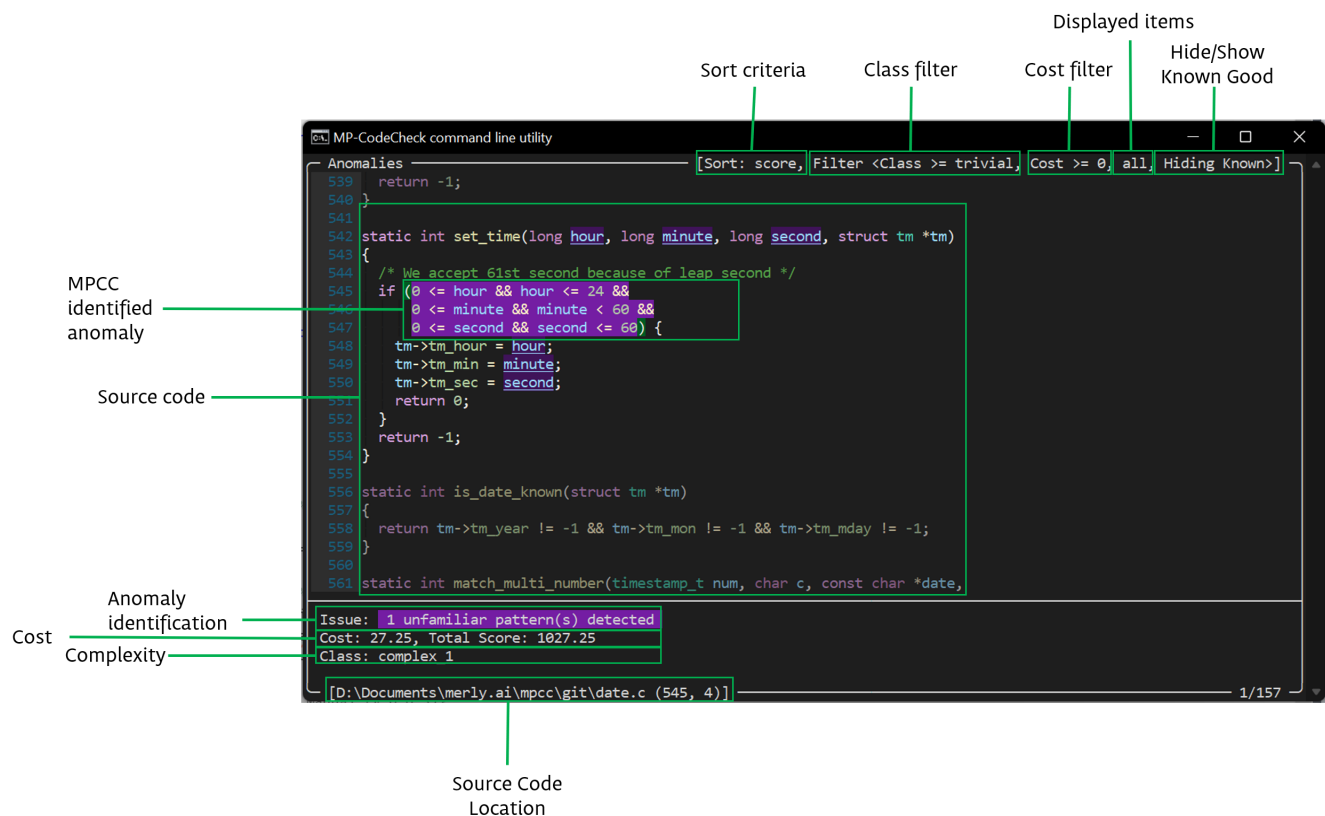


Figure 4. An Anomalous Example in MP-CodeCheck's Code View.

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.

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 10.

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.

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.

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

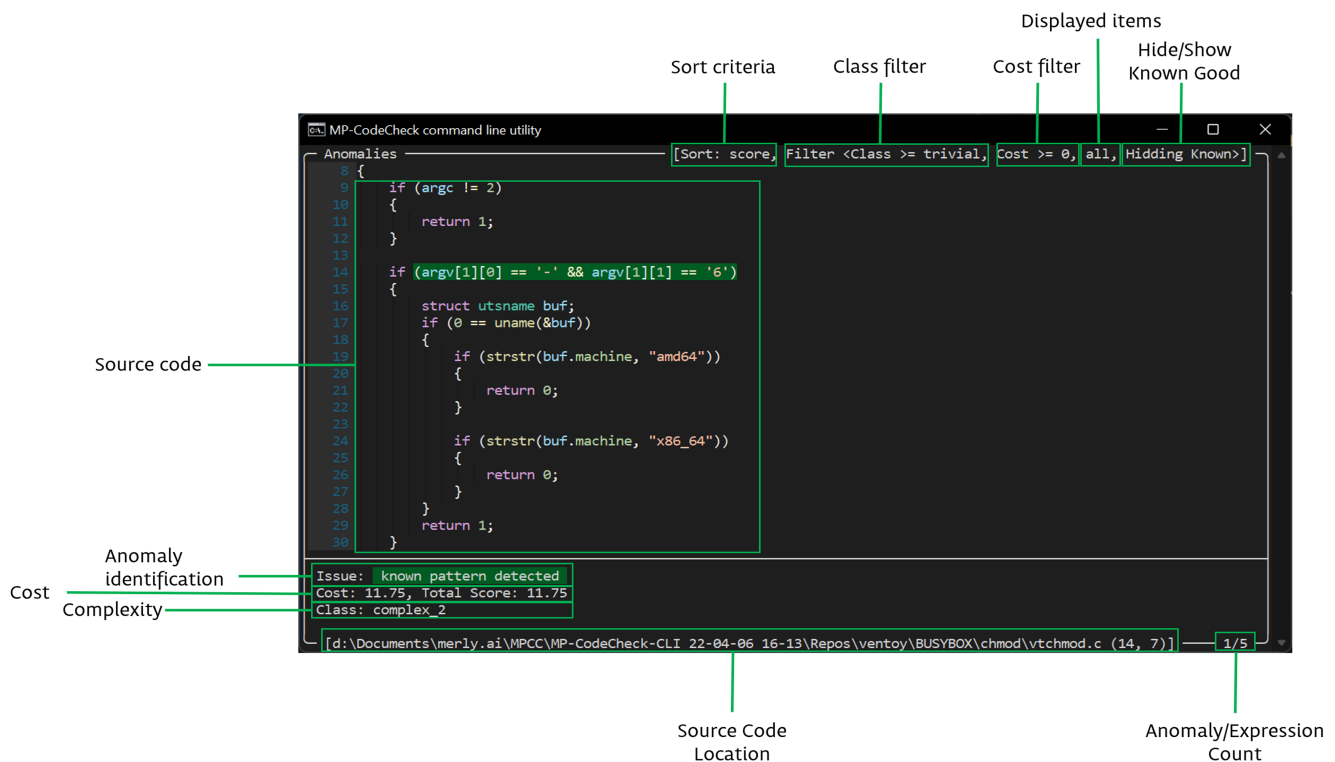


Figure 5. A Non-Anomalous Example in MP-CodeCheck’s Code View.

specific anomaly.

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.

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.

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.

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.

Sorting/Filtering Inference Results

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

```
MP-CodeCheck command line utility [Sort: score, Filter <Class >= trivial, Cost >= 0, all, Hiding Known>]
119 } else if (!strcmp(argv[arg], "--recover")) {
120     get_recover = 1;
121 } else if (!strcmp(argv[arg], "--stdin")) {
122     commits_on_stdin = 1;
123 } else if (skip_prefix(argv[arg], "--packfile=", &p)) {
124     const char *end;
125
126     packfile = 1;
127     if (parse_oid_hex(p, &packfile_hash, &end) || *end)
128         die_("argument to --packfile must be a valid hash (got '%s')", p);
129 } else if (skip_prefix(argv[arg], "--index-pack-arg=", &p)) {
130     strvec_push(&index_pack_args, p);
131 }
132 arg++;
133 }
134 if (argc != arg + 2 - (commits_on_stdin || packfile))
135     usage(http_fetch_usage);
136
137 if (nongit)
138     die_("not a git repository");
139
140 git_config(git_default_config, NULL);
141
142 if (packfile) {
143     if (!index_pack_args.nr)
144         die_("the option '%s' requires '%s'", "--packfile", "--index-pack-args");
145
146     fetch_single_packfile(&packfile_hash, argv[arg],
147                          index_pack_args.v);

```

Issue: known pattern detected
Cost: 9.25, Total Score: 9.25
Class: complex_2

[D:\Documents\merly.ai\MPCC\git\http-fetch.c (127, 6)] 2/24

Figure 6. MP-CodeCheck's Code View.

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

```

MP-CodeCheck command line utility - (42) D:\Documents\merly.ai\MPCC\git: 31664 Expressions / 42 Anomalies
Expr= (argc != arg + 2 - (commits_on_stdin || packfile)) 2008.95
1 | D:\Documents\merly.ai\MPCC\git\http-fetch.c (Line: 134)
-----
Expr= (len <= suffix_len || (p = git_dir + len - suffix_len)[-1] != '/' || strcmp(p, submodule_name)) 1060.95
2 | D:\Documents\merly.ai\MPCC\git\submodule.c (Line: 2225)
-----
Expr= (orig->size > MAX_XDIFF_SIZE || src1->size > MAX_XDIFF_SIZE || src2->size > MAX_XDIFF_SIZE || buffer_is_bin...1043.00
3 | D:\Documents\merly.ai\MPCC\git\ll-merge.c (Line: 111)
-----
Expr= ((p->token != GREP_PATTERN) && (p->token == GREP_PATTERN_HEAD == (ctx == GREP_CONTEXT_HEAD))) 1039.50
4 | D:\Documents\merly.ai\MPCC\git\grep.c (Line: 883)
-----
Expr= (last && last->data.len && last->data.buf && last->depth < max_depth && dat->len > the_hash_algo->rawsz) 1034.50
5 | D:\Documents\merly.ai\MPCC\git\builtin\fast-import.c (Line: 971)
-----
Expr= (reuse_delta && IN_PACK(trg_entry) && IN_PACK(trg_entry) == IN_PACK(src_entry) && !src_entry->preferred_bas...1031.75
6 | D:\Documents\merly.ai\MPCC\git\builtin\pack-objects.c (Line: 2477)
-----
Expr= (!strcmp(arg, "--all") || !strcmp(arg, "--branches") || !strcmp(arg, "--tags") || !strcmp(arg, "--remotes")...1029.75
7 | D:\Documents\merly.ai\MPCC\git\revision.c (Line: 2178)
-----
Expr= (0 <= hour && hour <= 24 && 0 <= minute && minute < 60 && 0 <= second && second <= 60) 1027.25
8 | D:\Documents\merly.ai\MPCC\git\date.c (Line: 545)
-----
Expr= (!delete && !rename && !copy && !edit_description && !new_upstream && !show_current && !unset_upstream && a...1026.25
9 | D:\Documents\merly.ai\MPCC\git\builtin\branch.c (Line: 723)
-----
Expr= (revs->min_age != -1 && (commit->date > revs->min_age) && !revs->line_level_traverse) 1024.50
10 | D:\Documents\merly.ai\MPCC\git\revision.c (Line: 1440)
-----
Expr= (!o->skip_sparse_checkout && (ce->ce_flags & CE_SKIP_WORKTREE) && (ce->ce_flags & CE_NEW_SKIP_WORKTREE)) 1023.25
11 | D:\Documents\merly.ai\MPCC\git\unpack-trees.c (Line: 2131)
-----
Expr= (revs.count && (revs.tag_objects || revs.tree_objects || revs.blob_objects) && (revs.left_right || revs.che...1020.75
[Sort: score, Filter <Class >= trivial, Cost >= 0, all, Hiding Known>]

```

Figure 7. MP-CodeCheck's Anomalies View.

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

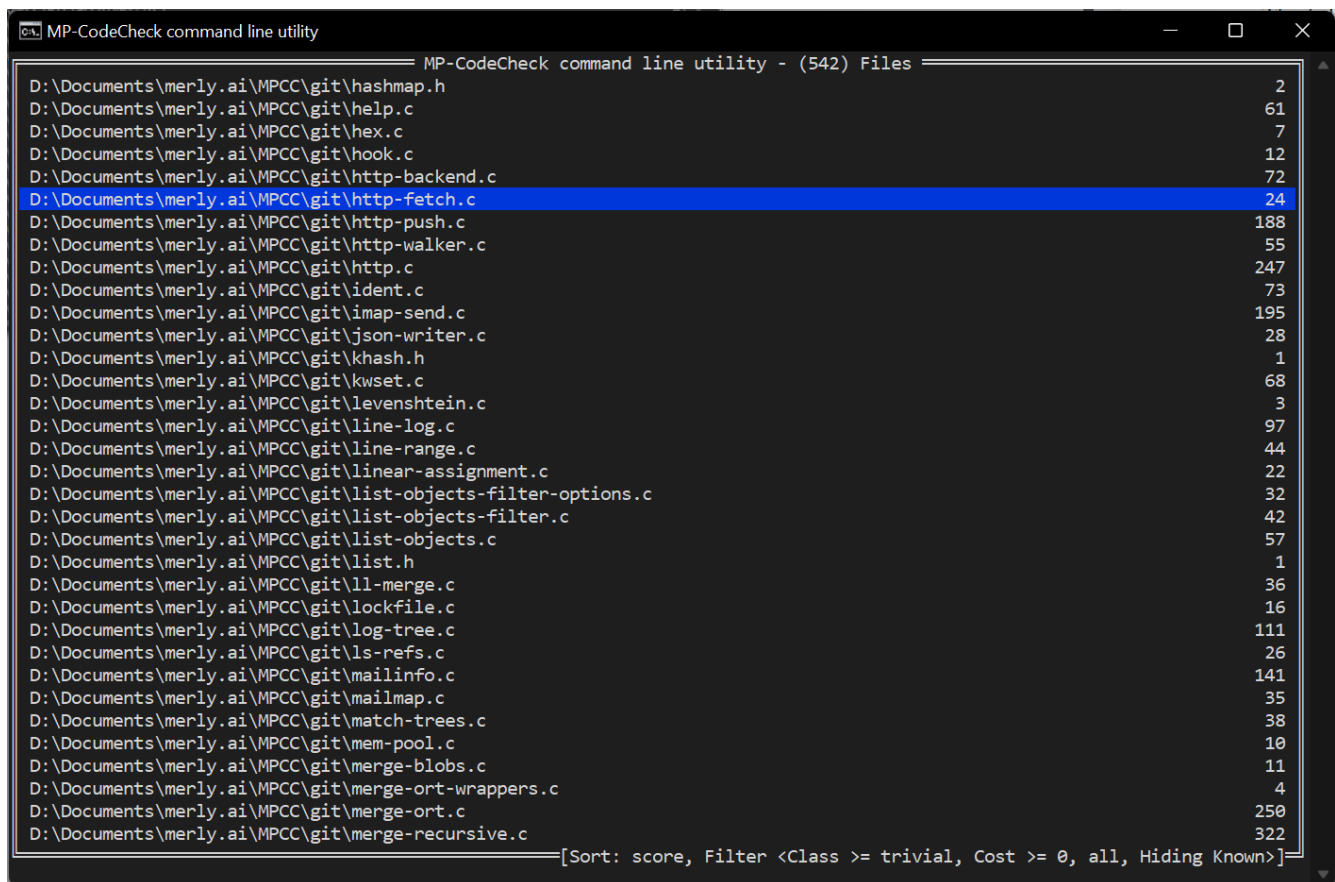


Figure 8. MP-CodeCheck's Files View.

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.

Expression	Score
(argc != arg + 2 - (commits_on_stdin packfile))	2008.95
(parse_oid_hex(p, &packfile_hash, &end) *end)	9.25
((ret = finish_http_pack_request(preq)))	7.50
(argv[arg][1] == 't')	5.75
(argv[arg][1] == 'c')	5.75
(argv[arg][1] == 'a')	5.75
(argv[arg][1] == 'v')	5.75
(argv[arg][1] == 'w')	5.75
(argv[arg][1] == 'h')	5.75
(skip_prefix(argv[arg], "--packfile=", &p))	4.25
(skip_prefix(argv[arg], "--index-pack-arg=", &p))	4.25
(!nurl !git_env_bool("GIT_TRACE_REDACT", 1))	4.00
(results.curl_result != CURLE_OK)	3.75
(start_active_slot(preq->slot))	3.50
(walker->corrupt_object_found)	3.25
(!strcmp(argv[arg], "--recover"))	2.75
(!strcmp(argv[arg], "--stdin"))	2.75
(preq == NULL)	2.25
(!index_pack_args.nr)	2.25
(index_pack_args.nr)	1.50
(commits_on_stdin)	0.00
(nongit)	0.00
(packfile)	0.00
(commits_on_stdin)	0.00

Figure 9. MP-CodeCheck’s Expressions View.

[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.

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.

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

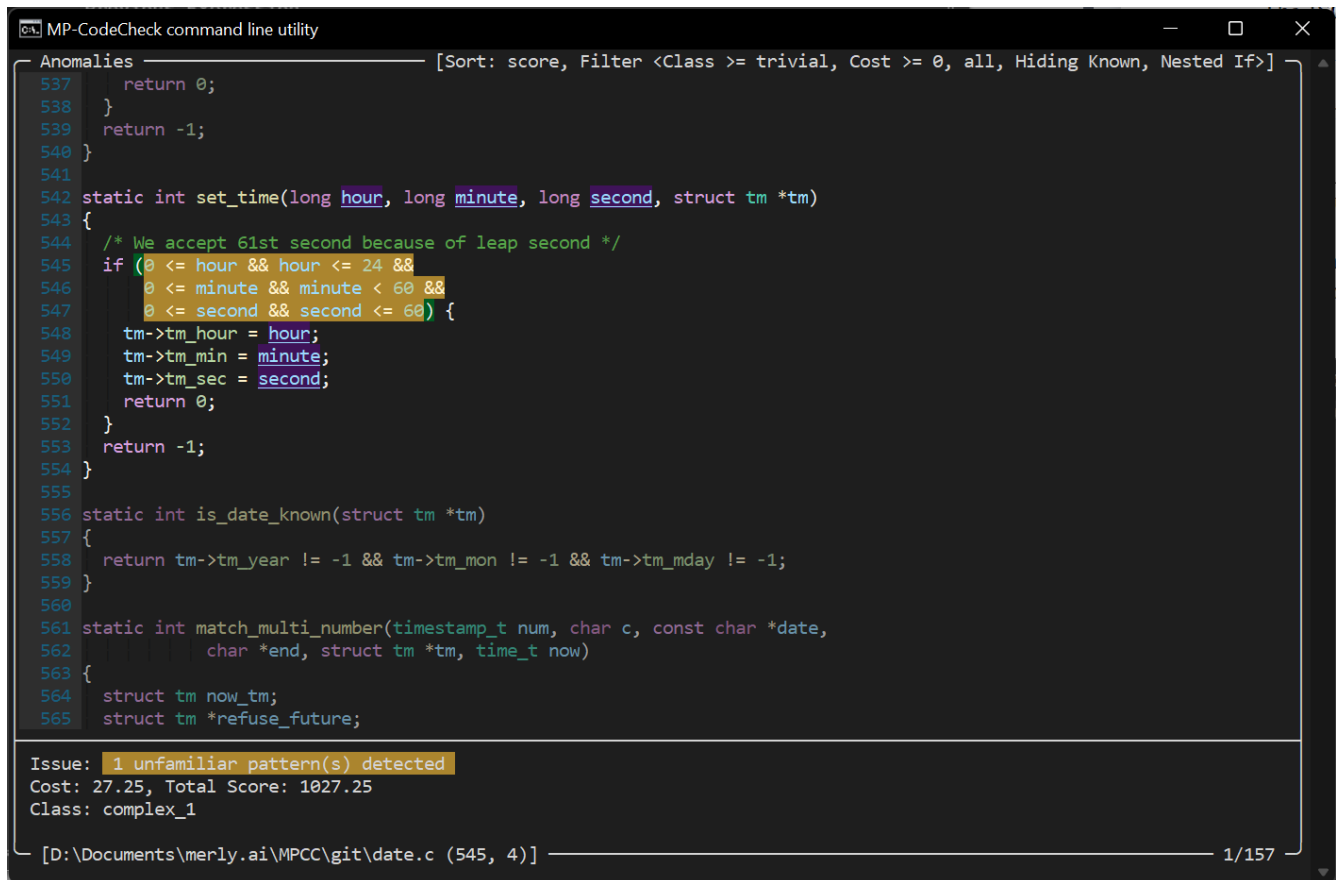
Log Files: You can change the model path by setting the directory associated with: `local-db_root_path`



Figure 10. MP-CodeCheck's Help Dialogue Box.

You can change the log path by setting the directory associated with: log_path

Settings: “run-training” – Determines whether or not training should be run before inference on the source code (defaults to true). “filter” – Determines whether or not the nested.ifs are extracted (defaults to true).



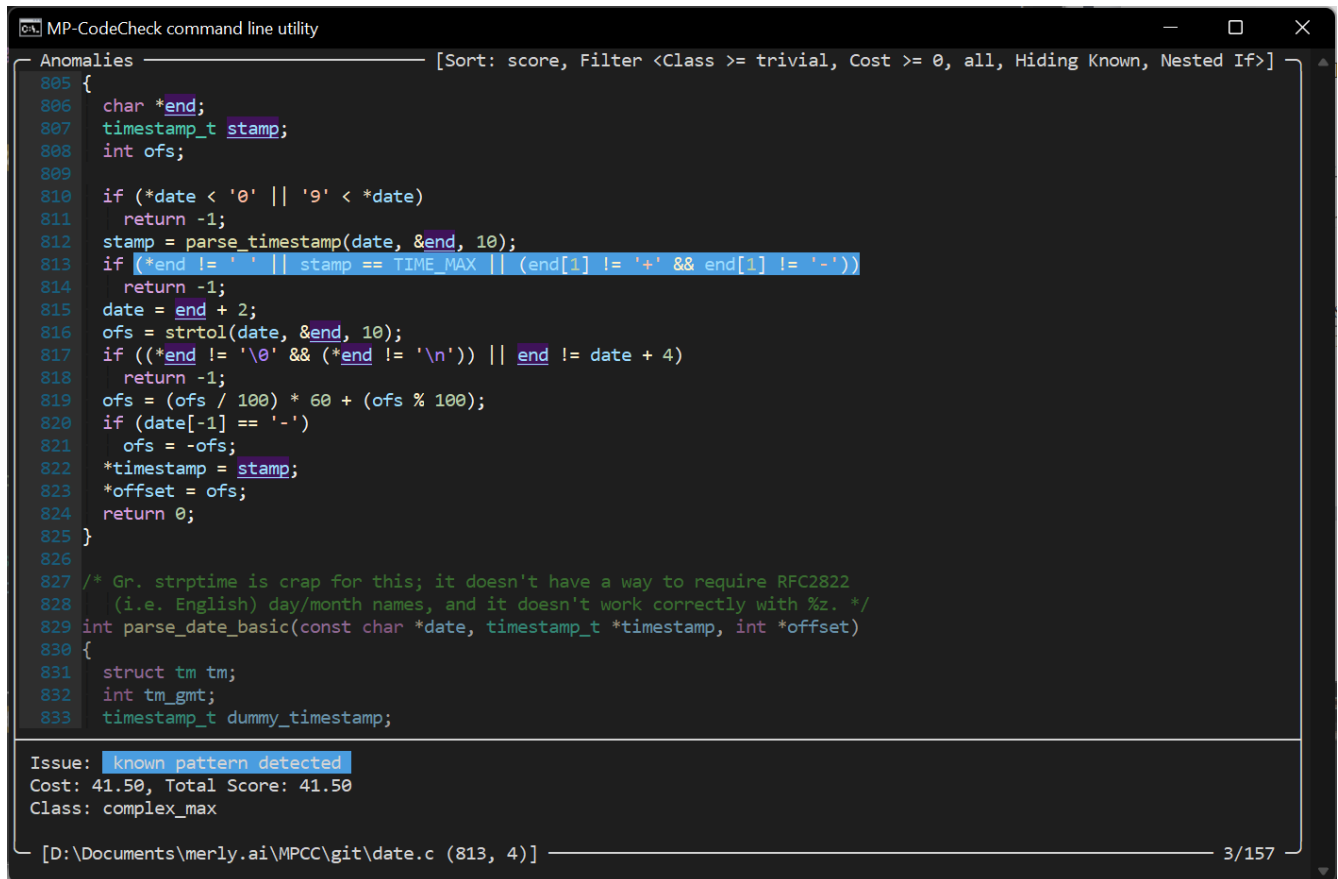
The screenshot shows a window titled "MP-CodeCheck command line utility". At the top, there is a tab labeled "Anomalies" and a filter bar with the text "[Sort: score, Filter <Class >= trivial, Cost >= 0, all, Hiding Known, Nested If>]". The main area displays C code from a file named "date.c". The code includes functions like `set_time`, `is_date_known`, and `match_multi_number`. An anomaly is highlighted in orange on line 545, which is the opening curly brace of an `if` statement. The bottom of the window shows a summary of the issue: "Issue: 1 unfamiliar pattern(s) detected", "Cost: 27.25, Total Score: 1027.25", and "Class: complex_1". The status bar at the bottom indicates the file path "D:\Documents\merly.ai\MPCC\git\date.c (545, 4)" and the page number "1/157".

```
537     return 0;
538 }
539 return -1;
540 }
541
542 static int set_time(long hour, long minute, long second, struct tm *tm)
543 {
544     /* We accept 61st second because of leap second */
545     if (0 <= hour && hour <= 24 &&
546         0 <= minute && minute < 60 &&
547         0 <= second && second <= 60) {
548         tm->tm_hour = hour;
549         tm->tm_min = minute;
550         tm->tm_sec = second;
551         return 0;
552     }
553     return -1;
554 }
555
556 static int is_date_known(struct tm *tm)
557 {
558     return tm->tm_year != -1 && tm->tm_mon != -1 && tm->tm_mday != -1;
559 }
560
561 static int match_multi_number(timestamp_t num, char c, const char *date,
562                               char *end, struct tm *tm, time_t now)
563 {
564     struct tm now_tm;
565     struct tm *refuse_future;
```

Issue: 1 unfamiliar pattern(s) detected
Cost: 27.25, Total Score: 1027.25
Class: complex_1

[D:\Documents\merly.ai\MPCC\git\date.c (545, 4)] 1/157

Figure 11. An example of changing the anomaly highlight color in MP-CodeCheck's Code View.



The screenshot shows a window titled "MP-CodeCheck command line utility". At the top, there is a filter bar with the text "[Sort: score, Filter <Class >= trivial, Cost >= 0, all, Hiding Known, Nested If]". The main area displays C code from lines 805 to 833. Line 813 is highlighted in blue, indicating an issue. The code is as follows:

```
805 {
806     char *end;
807     timestamp_t stamp;
808     int ofs;
809
810     if (*date < '0' || '9' < *date)
811         return -1;
812     stamp = parse_timestamp(date, &end, 10);
813     if (*end != ' ' || stamp == TIME_MAX || (end[1] != '+' && end[1] != '-'))
814         return -1;
815     date = end + 2;
816     ofs = strtol(date, &end, 10);
817     if ((*end != '\0' && (*end != '\n')) || end != date + 4)
818         return -1;
819     ofs = (ofs / 100) * 60 + (ofs % 100);
820     if (date[-1] == '-')
821         ofs = -ofs;
822     *timestamp = stamp;
823     *offset = ofs;
824     return 0;
825 }
826
827 /* Gr. strptime is crap for this; it doesn't have a way to require RFC2822
828    (i.e. English) day/month names, and it doesn't work correctly with %z. */
829 int parse_date_basic(const char *date, timestamp_t *timestamp, int *offset)
830 {
831     struct tm tm;
832     int tm_gmt;
833     timestamp_t dummy_timestamp;
```

Below the code, the following information is displayed:

Issue: known pattern detected
Cost: 41.50, Total Score: 41.50
Class: complex_max

At the bottom of the window, the file path and line information are shown: "[D:\Documents\merly.ai\MPCC\git\date.c (813, 4)]" and "3/157".

Figure 12. An example of changing the expression highlight color in MP-CodeCheck’s Code View.



THE MACHINE PROGRAMMING COMPANY