

CPSC 457 - Assignment 2

Due date is posted on D2L.

Individual assignment. Group work is NOT allowed.

Weight: 19% of the final grade.

Overall description

For this assignment you will write a C++ function `getDirStats()` that recursively examines the contents of a directory and then reports some statistics about the files and subdirectories it discovers, such as largest file and duplicate files. You will be given an incomplete C++ program that contains an unfinished `getDirStats()` function, and other helpful functions and classes. Your job is to finish the program by implementing the `getDirStats()` function. This function will compute: (a) largest file in the directory, (b) size of the largest file, (c) total number of files in the directory, (d) total number of subdirectories, (e) sum of all file sizes, (f) list of 5 most common file types, (g) list of 5 largest groups of duplicate files.

Details

Start by cloning the code in the GitLab repository:

```
$ git clone https://gitlab.com/cpsc457/public/dirstats.git
```

Then compile the code using included `Makefile`, to create an executable `dirstats`.

```
$ cd dirstats
$ make
```

Then run the executable.

```
$ ./dirstats .
-----
Largest file:      ./some_dir/some_file.txt
Largest file size: 123
Number of files:   321
Number of dirs:    333
Total file size:   1000000
Most common file types:
- C source, ASCII text
- makefile script, ASCII text
- C++ source, ASCII text
- directory
Duplicate files - group 1:
- ./file1.cpp
- ./lib/sub/other.c
Duplicate files - group 2:
- ./readme.md
- ./docs/readme.txt
- ./x.y
```

The executable is created from a `main()` function defined in `main.cpp`, which parses the command line arguments, calls an unfinished `getDirStats()` function, and then formats the results on standard output.

You can only modify the `getDirStats.cpp` file, and you must not modify any other files.

The `getDirStats()` function

The `getDirStats()` is declared in `getDirStats.h` and defined in `getDirStats.cpp`:

```
bool getDirStats(const std::string & dir_name, Results & res);
```

It takes two parameters:

- `dir_name` is the top directory that your function must examine recursively; and
- `res` is a reference to a structure where your function will store results.

The function returns:

- `true` on success, which means the results will be available to the caller in the `res` parameter;
- `false` if the function could not finish, e.g. because some of the files/directories were unreadable. In this case the `res` parameter will contain undefined values.

The `res` parameter is of type `struct Results`, defined in `getDirStats.h`, which will allow your function to return multiple results, as described below:

- `std::string largest_file_path`
In this field you will return a path to the largest file in the directory. The path needs to include the top directory name. If the directory contains no files, the value should be set to an empty string.
- `long largest_file_size`
You will use this field to return the size of the largest file in the directory. If directory contains no files, the value should be set to -1.
- `long n_files`
Will contain the total number of files encountered while scanning the directory.
- `long n_dirs`
Will contain the total number of sub-directories encountered while scanning the directory.
- `long all_files_size`
Will contain the sum of all sizes of all files encountered during the scan. Use `stat()` to determine file sizes for each file, then sum them up.
- `std::vector<std::string> most_common_types`
Will contain a list of 5 most common file type encountered, sorted in descending order by the frequency of the type. See below for more details.
- `std::vector<std::vector<std::string>> duplicate_files`
Will contain a list of 5 largest groups of duplicate files, sorted in descending order by the size of the groups. See below for more details.

Computing the first 5 fields

These are very simple to compute. You need to write code that traverses the directory recursively, and calls the `stat()` function on all directory entries encountered. You then extract some information from the returned structure and update the results. You can use the code from <https://gitlab.com/cpsc457/public/find-empty-directories> as a starting point.

The remaining 2 fields are not too difficult either, but they require a bit more explanation.

Computing `most_common_types`

You will need to use `popen()` to call UNIX `file(1)` utility with the `-b` option to obtain the file type of the file. Your code will need to parse the output of `file(1)` and only consider the main type (the text before the first comma). Once you have the file types for all files, you will need to find the 5 most common types. There are many ways to do this, but I suggest you use `std::unordered_map<std::string,int>` container and `std::sort()` if you want to make your life easy.

Computing `duplicate_files`

You could try to compare the contents of every file to every other file, but this would be slow, $O(n^2)$!!! A better way, $O(n \log n)$, is to compute a hash of every file, sort the hashes, and then discover duplicates by detecting the same consecutive hashes. To get a hash of a file you can use the provided function:

```
std::string sha256_from_file( const std::string & fname);
```

This function computes a SHA256 digest for a given file and returns it as a string. It mimics some of the functionality of the `sha256sum` UNIX utility, described in little more detail in the appendix.

Another efficient way to find duplicates would be to use a dictionary that maps digests to an array of filenames, e.g. using:

```
std::unordered_map<std::string,std::vector<std::string>>
```

The `duplicate_files` field will contain the list of 5 largest groups of identical files. Each group will be represented as a list of files (`std::vector<std::string>`) belonging to that group. No matter which technique you use to identify the groups of identical files, eventually you will need to sort these groups by size (e.g. using `std::sort`). Then you populate the `duplicate_files` with the largest groups such that:

- `duplicate_files.size()` will return the number of groups of duplicate files, at most 5
- `duplicate_files[0]` will contain the largest group of identical files
- `duplicate_files[i].size()` will return the number of duplicate files in the *i*-th group
- `duplicate_files[i][j]` will return the filepath of the *j*-th file in the *i*-th group of duplicates

For example, to return 2 groups of files ["a", "b", "c"] and ["x", "y"], you could write code like this:

```
std::vector<std::string> group1;
group1.push_back("a"); group1.push_back("b"); group1.push_back("c");
res.duplicate_files.push_back(group1);
std::vector<std::string> group2;
group2.push_back("x"); group2.push_back("y");
res.duplicate_files.push_back(group2);
```

Important

All code you write must go in the `getDirStats.cpp` file, and that should be the only file you will submit for grading. Your TAs will test your code by supplying their own `main()` function, which may be different from the `main()` that you will be using. It is therefore vital that you maintain the same function signature as declared in `getDirStats.h`. Before you submit `getDirStats.cpp` to D2L, make sure it works with the provided `main()` function!!!

Assumptions:

- None of the files and directory names will contain spaces.
- The total number of directories and files will be less than 5000.
- Each filename will contain less than 4KiB characters.

Grading

Your code will be graded on correctness and efficiency. Your code should be at least as efficient as the included Python solution (described in the appendix).

Allowed libraries

You are free to use any APIs from the lib/libc++ libraries for this assignment. This includes, but is not limited to:

- `popen()` to get the output of the `file` utility
- `stat()`, `opendir()`, `closedir()`, `readdir()`, `getcwd()`, `chdir()`
- `open()`, `close()`, `read()`, `fopen()`, `fread()`, `fclose()`
- `std::map`, `std::unordered_map`, `std::vector`, `std::string`, `std::sort`
- C++ streams

Submission

Submit `getDirStats.cpp` file to D2L.

Appendix 1 – computing file digests using `sha256sum(1)` utility

Cryptographic hash functions are special hash functions that have many different uses. You will be using them, namely SHA256, to efficiently detect duplicate file contents. If you have never experimented with file digests, there is a UNIX utility `sha256sum` that allows you to compute digests for files from the command line. Here is an example of how to use it from command line to get a digest of a file `1.txt`. Let's start by creating a sample file `1.txt` and computing its digest:

```
$ echo "hello" > 1.txt
$ sha256sum 1.txt
5891b5b522d5df086d0ff0b110fbd9d21bb4fc7163af34d08286a2e846f6be03 1.txt
```

The long **hexadecimal** string is the digest of the file contents. The idea behind the digest is that it will be different for every file – as long as the file contents are different. Here is an example of a digest for a slightly modified file:

```
$ echo "world" >> 1.txt
$ sha256sum 1.txt
4a1e67f2fe1d1cc7b31d0ca2ec441da4778203a036a77da10344c85e24ff0f92 1.txt
```

As you can see, the digest is now completely different. We can use digests to compare two really big files – all we have to do is compute the digests for both files, and if the digests are different, then the file contents are definitely different. If the digests are the same, there is a very high probability that the files are the same. What is the likelihood of getting the same digests for different file contents? It is very small. So small that git uses a weaker, 128-bit version of SHA, to detect differences in files.

How does this help us detect duplicate files? Instead of comparing the contents of every single file to every other file, which would be slow, we can instead compute the digests for all files first, then sort the digests + filename pairs by the digests, and then detect duplicates by scanning the result once. You can even do this on the command line. For example, to find all duplicates in directory `/usr/include/c++` you could execute:

```
$ find /usr/include/c++ -type f | xargs sha256sum \
  | sort | uniq -w64 --all-repeated=separate
...
acaa059f12b827ac1b071d9b7ce40c7043e4e65c5b4ca2169752cd81d8993356
/usr/include/c++/9/backward/hash_set
acaa059f12b827ac1b071d9b7ce40c7043e4e65c5b4ca2169752cd81d8993356
/usr/include/c++/9/ext/hash_set

c0b15be47a20948066531a4570f13fe02dbe0f7580cd8b5cad5fb172fabe9f33
/usr/include/c++/9/x86_64-redhat-linux/32/bits/gthr-default.h
c0b15be47a20948066531a4570f13fe02dbe0f7580cd8b5cad5fb172fabe9f33
/usr/include/c++/9/x86_64-redhat-linux/32/bits/gthr-posix.h
c0b15be47a20948066531a4570f13fe02dbe0f7580cd8b5cad5fb172fabe9f33
/usr/include/c++/9/x86_64-redhat-linux/bits/gthr-default.h
c0b15be47a20948066531a4570f13fe02dbe0f7580cd8b5cad5fb172fabe9f33
/usr/include/c++/9/x86_64-redhat-linux/bits/gthr-posix.h
...
```

The output shows all duplicate files in groups. For example, you can see that the file `gthr-default.h` is repeated verbatim 4 times.

Appendix 2 – obtaining/guessing file types using `file(1)` utility

The `file(1)` utility is used on UNIX systems to guess the file type. It prints out the main type of the file, optionally followed by extra details, separated by commas. For example, to see all types of the files in the git repo you could type:

```
$ file *
digester.cpp:      C source, ASCII text
digester.h:        C++ source, ASCII text
getDirStats.cpp:   C source, ASCII text
getDirStats.h:     C source, ASCII text
main.cpp:          C source, ASCII text
Makefile:          makefile script, ASCII text
README.md:         ASCII text
test1:             directory
```

The output above identified `digester.h` file as a “C++ source”, with the extra detail “ASCII text”. For this assignment we only care about the main type, i.e. “C++ source”. To make parsing easier, you can invoke the file utility with the `-b` option, which will omit the filenames from output.

Here is an example of a command line you can use to get the most common types from a large directory `/usr/share/doc/octave`:

```
$ find /usr/share/doc/octave/ -type f \
| xargs file -b \
| awk -F, '{print $1}' \
| sort \
| uniq -c \
| sort -nr

2651 HTML document
  36 C source
  29 ASCII text
  27 PNG image data
   5 PDF document
   3 UTF-8 Unicode text
   1 LaTeX document
   1 C++ source
```

The first `find` command in the pipe recursively scans the directory and prints all files. The second `xargs` command runs ‘`file -b`’ on every file. The third `awk` command splits the output by commas and prints only the first field. The fourth `sort` command sorts the output alphabetically. The fifth `uniq` command counts consecutive line occurrences and prints the counts next to the strings. The last command `sort` sorts the output numerically, in descending order.

Appendix 3 – python solution

The repository includes a Python program `dirstats.py` that solves the assignment. This should help you design your own test cases and see what the expected output should look like. Here is an example of running it on the test1 directory:

```
$ ./dirstats.py test1
Largest file:      test1/dir1/down.png
Largest file size: 202
Number of files:   7
Number of dirs:    9
Total file size:   253
Most common file types:
- ASCII text
- PNG image data
- empty
Duplicate file - group 1
- test1/dir1/file1.txt
- test1/dir1/file2.txt
- test1/a/b/c/d/f
```

General information about all assignments:

1. All assignments are due on the date listed on D2L. Late submissions will be not be marked.
2. Extensions may be granted only by the course instructor.
3. After you submit your work to D2L, verify your submission by re-downloading it.
4. You can submit many times before the due date. D2L will simply overwrite previous submissions with newer ones. It is better to submit incomplete work for a chance of getting partial marks, than not to submit anything. Please bear in mind that you cannot re-submit a single file if you have already submitted other files. Your new submission would delete the previous files you submitted. So please keep a copy of all files you intend to submit and resubmit all of them every time.
5. Assignments will be marked by your TAs. If you have questions about assignment marking, contact your TA first. If you still have questions after you have talked to your TA, then you can contact your instructor.
6. All programs you submit must run on linux.cpsc.ucalgary.ca. If your TA is unable to run your code on the Linux machines, you will receive 0 marks for the relevant question.
7. Unless specified otherwise, you must submit code that can finish on any valid input under 10s on linux.cpsc.ucalgary.ca, when compiled with `-O2` optimization. Any code that runs longer than this may receive a deduction, and code that runs for too long (about 30s) will receive 0 marks.
8. **Assignments must reflect individual work.** For further information on plagiarism, cheating and other academic misconduct, check the information at this link: <http://www.ucalgary.ca/pubs/calendar/current/k-5.html>.
9. Here are some examples of what you are not allowed to do for individual assignments: you are not allowed to copy code or written answers (in part, or in whole) from anyone else; you are not allowed to collaborate with anyone; you are not allowed to share your solutions with anyone else; you are not allowed to sell or purchase a solution. This list is not exclusive.
10. We will use automated similarity detection software to check for plagiarism. Your submission will be compared to other students (current and previous), as well as to any known online sources. Any cases of detected plagiarism or any other academic misconduct will be investigated and reported.