

609. Find Duplicate File in System

DescriptionHintsSubmissionsSolutions

- Total Accepted: **2198**
- Total Submissions: **4162**
- Difficulty: **Medium**
- Contributors:fallcreek

Given a list of directory info including directory path, and all the files with contents in this directory, you need to find out all the groups of duplicate files in the file system in terms of their paths.

A group of duplicate files consists of at least **two** files that have exactly the same content.

A single directory info string in the **input** list has the following format:

```
"root/d1/d2/.../dm f1.txt(f1_content) f2.txt(f2_content) ... fn.txt(fn_content)"
```

It means there are **n** files (**f1.txt**, **f2.txt** ... **fn.txt** with

content **f1_content**, **f2_content** ... **fn_content**, respectively) in directory **root/d1/d2/.../dm**. Note

that $n \geq 1$ and $m \geq 0$. If $m = 0$, it means the directory is just the root directory.

The **output** is a list of group of duplicate file paths. For each group, it contains all the file paths of the files that have the same content. A file path is a string that has the following format:

```
"directory_path/file_name.txt"
```

Example 1:

Input:

```
["root/a 1.txt(abcd) 2.txt(efgh)", "root/c 3.txt(abcd)", "root/c/d 4.txt(efgh)", "root 4.txt(efgh)"]
```

Output:

```
[["root/a/2.txt","root/c/d/4.txt","root/4.txt"],["root/a/1.txt","root/c/3.txt"]]
```

Note:

1. No order is required for the final output.
2. You may assume the directory name, file name and file content only has letters and digits, and the length of file content is in the range of [1,50].
3. The number of files given is in the range of [1,20000].
4. You may assume no files or directories share the same name in the same directory.
5. You may assume each given directory info represents a unique directory. Directory path and file info are separated by a single blank space.

Follow-up beyond contest:

1. Imagine you are given a real file system, how will you search files? DFS or BFS?
2. If the file content is very large (GB level), how will you modify your solution?
3. If you can only read the file by 1kb each time, how will you modify your solution?
4. What is the time complexity of your modified solution? What is the most time-consuming part and memory consuming part of it? How to optimize?
5. How to make sure the duplicated files you find are not false positive?

[Subscribe](#) to see which companies asked this question.

```
vector<string> strsplit(string path,string tok)
{
    vector<string> res;
    char *ss = strtok((char*)path.c_str(),(char *)tok.c_str());
    res.push_back(string(ss));

    while(ss!=NULL)
    {
        ss = strtok(NULL,(char *)tok.c_str());
        if(ss!=NULL)
            res.push_back(string(ss));
    }
}
```

```

    }
    return res;
}

vector<vector<string>> findDuplicate(vector<string>& paths) {
    unordered_map<string,vector<string>> pathMap;
    for(auto path:paths)
    {
        vector<string> splitarr = strsplit(path," ");
        for(int i = 1;i<(int)splitarr.size();i++)
        {
            string pathtmp = splitarr[i];
            int startidx = pathtmp.find("(");
            int endidx = pathtmp.find(")");
            string content = pathtmp.substr(startidx+1,endidx-startidx-1);
            string subpath = pathtmp.substr(0,startidx);

            //cout<<subpath<<endl;
            string pathstr = splitarr[0]+"/"+subpath;
            pathMap[content].push_back(pathstr);
        }
        //pathMap[substr].push_back(content);
    }
    vector<vector<string>> res;
    for(auto iter:pathMap)
    {
        vector<string> ans_tmp = iter.second;
        if(ans_tmp.size()>1)
            res.push_back(ans_tmp);
    }
    return res;
}

```