# Practice 9 (2015/12/1)
## SPAM email checker

1. **File IO**

   A. ***fopen, fclose, fgets***

   B. **Example**

```
1   #include <stdio.h>
2   #define MAX_LEN 200
3   #define MAX_ROW 100
4
5   int main(int argc, char **argv)
6   {
7       FILE * pFile;
8       char mystring [MAX_ROW][MAX_LEN];
9
10      pFile = fopen (*(argv+1) , "r");
11      if (pFile == NULL) printf("Error opening file\n");
12      else {
13        int i = 0;
14        while ( fgets (mystring[i] , MAX_LEN-1 , pFile) != NULL )  {
15          puts (mystring[i]); i++;
16        }
17        fclose (pFile);
18      }
19      return 0;
20  }
```

```
After the last annual calculations of your account activity we have determined
that you are eligible to receive a tax refund of $479.30 . Please submit the
tax refund request and allow us 2-6 days in order to process it.

A refund can be delayed for a variety of reasons.
For example submitting invalid records or applying after the deadline.
To access the form for your tax refund, please click here
(http://e-dlogs.rta.mi.th:84/www.irs.gov/)

Note: Deliberate wrong inputs will be prosecuted by law.

Regards,
Internal Revenue Service
```

After the last annual calculations of your account activity we have determined

that you are eligible to receive a tax refund of $479.30 . Please submit the

tax refund request and allow us 2-6 days in order to process it.


A refund can be delayed for a variety of reasons.

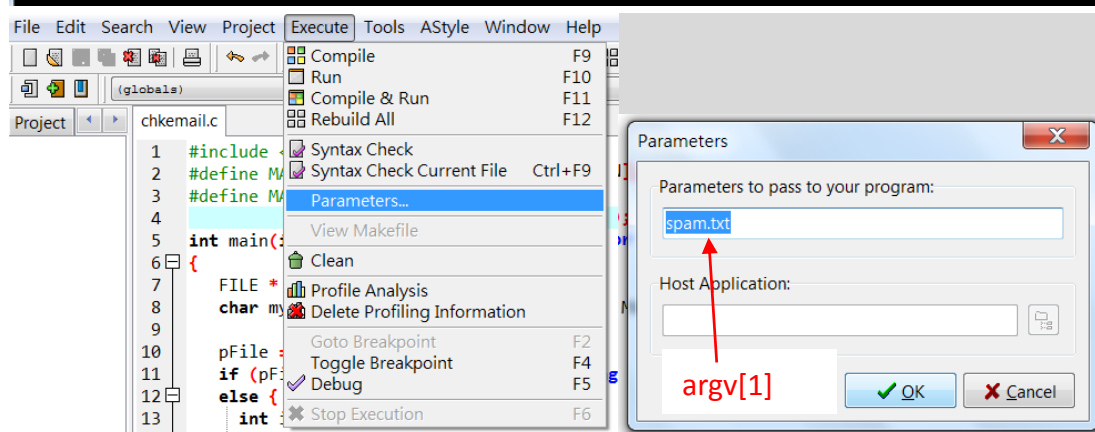For example submitting invalid records or applying after the deadline.

To access the form for your tax refund, please click here

(http://e-dlogs.rta.mi.th:84/www.irs.gov/)


Note: Deliberate wrong inputs will be prosecuted by law.


Regards,

Internal Revenue Service

--------------------------------

File  Edit  Search  View  Project  Execute  Tools  AStyle  Window  Help

Compile          F9
Run              F10
Compile & Run    F11
Rebuild All      F12

Syntax Check
Syntax Check Current File    Ctrl+F9
Parameters...
View Makefile
Clean
Profile Analysis
Delete Profiling Information
Goto Breakpoint          F2
Toggle Breakpoint        F4
Debug                    F5
Stop Execution           F6

Project

chkemail.c

```
1   #include
2   #define M
3   #define M
4
5   int main(
6 { 
7       FILE *
8       char m
9
10      pFile
11      if (pF
12 {    else {
13          int
```

Parameters

Parameters to pass to your program:

spam.txt

Host Application:

argv[1]

OK    Cancel

## 2. Problem

SPAM emails consume many network resources and cost U.S. organization billions of dollars a year in spam-prevention software, equipment, network resources, bandwidth and lost productivity. Write a program to read an email and save its content into a 2-D character array (one line from the first character to the newline in a row). **Initially you are given several spam email examples (files), a set of keywords for each type of spam emails (each type in a file), and a set of weights for each spam keyword.** For each spam keyword, find out the number of occurrences in an email. Then you can calculate the total cost for an email and determine if this is a spam email. Here is a simple and straightforward way to define the weight for each keyword and compute the total weight for an email. For each keyword's weight, you can just simply use the appearing rate of a keyword

among all keywords as its weight. For instance, you have eleven keywords and all keywords appear 100 times in total in your spam email examples. For keyword A, it only appears 3 times, then A's weight is 0.03. As for the way to compute the spam weight for an email is to use this formula:

$$sw_e = \sum_{k \in KW} w_k \times \frac{a_k}{N_w}$$

Where $sw_e$ is the total spam weight for email $e$, $KW$ is the set of keywords, $w_k$ is the weight of keyword $k$, $N_w$ is **the total number of words of the email** under analysis, and $a_k$ is **the number of occurrences of keyword $k$ in this email**. As $sw_e$ is higher than a threshold value, this email is regarded a spam email.

```c
1   #include <stdio.h>
2   #include <string.h>
3   #define MAX_LEN 200
4   #define MAX_ROW 200
5
6   int main(int argc, char **argv)
7   {
8       FILE * pFile;
9       char keyword [MAX_ROW][MAX_LEN];
10
11      if ((pFile = fopen (*(argv+1) , "r")) == NULL) printf("Error opening file\n");
12      else {
13        int i = 0;
14        while ( 1 )    {
15          fscanf (pFile, "%s", keyword[i]);
16          puts (keyword[i]);
17          if (!strcmp(keyword[i],"-999"))  {
18            i--; break;
19          }
20          i++;
21        }
22        fclose (pFile);
23      }
24      return 0;
25  }
```

```
fund
email
exceed
quota
limit
account
-999

------------------------------
Process exited after 0.01082 seconds with return value 0
請按任意鍵繼續 . . . ■
```

**Main features:**

(a) **Your program has several features: (1) compute the number of occurrences for each keyword in an email and report the statistics; (2) report your judgement.**

3

**(b)** For feature (1), the program needs to read in a keyword file and an email file. Run your program like: ( due tonight)

*chkspam –k keyword_file_name –e email_file_name –oc*

where –k means the following parameter specifies a keyword file and –n means the following parameter specifies an email file.

```
***The statistics for the occurrences of each keyword in the file spam1.txt ***
    quota:  0
     size: 10
    limit:  3
    email:  1
******************************************************************************

-------------------------------
Process exited after 1.452 seconds with return value 0
請按任意鍵繼續 . . .
```

**(c)** Enhance feature (1) as (due in 12/7):

*chkspam –kn kFile_1 kFile_2 … kFile_n –em eFile_1 eFile_2 … eFile_m -oc*

where –kn means the following *n* parameters are *n* keyword files, -em means the following *m* parameters are *m* email files and –oc means the execution of reporting occurrence statistics. This feature can read many keyword and email files sequentially, merge all keywords as a single set, and identify the occurrences of each keyword in all email files. If you fail in opening any file, you have to prompt the failure message but the program must continue to go.

**(d)** For feature (2), use another option to check if an email file is a spam email. (due in 12/7)

*chkspam –kn kFile_1 kFile_2 … kFile_n -em eFile_1 … eFile_m –spam*

You can read one or several keyword files to one or several emails. You have to determine a threshold value as spam lower bound, which implies that if the total cost for an email exceeds that threshold value then that email is regarded as a spam email. You may also need to adjust the weight for some keywords to improve the correctness of your judgement.

**(e)** Enhance your keyword library and derive another new way, such as the use of sub-sentence checking, to improve the feasibility and correctness of your program to identify spam emails. (due in 12/7)