

Information Retrieval Assignment

Teaching Objective:

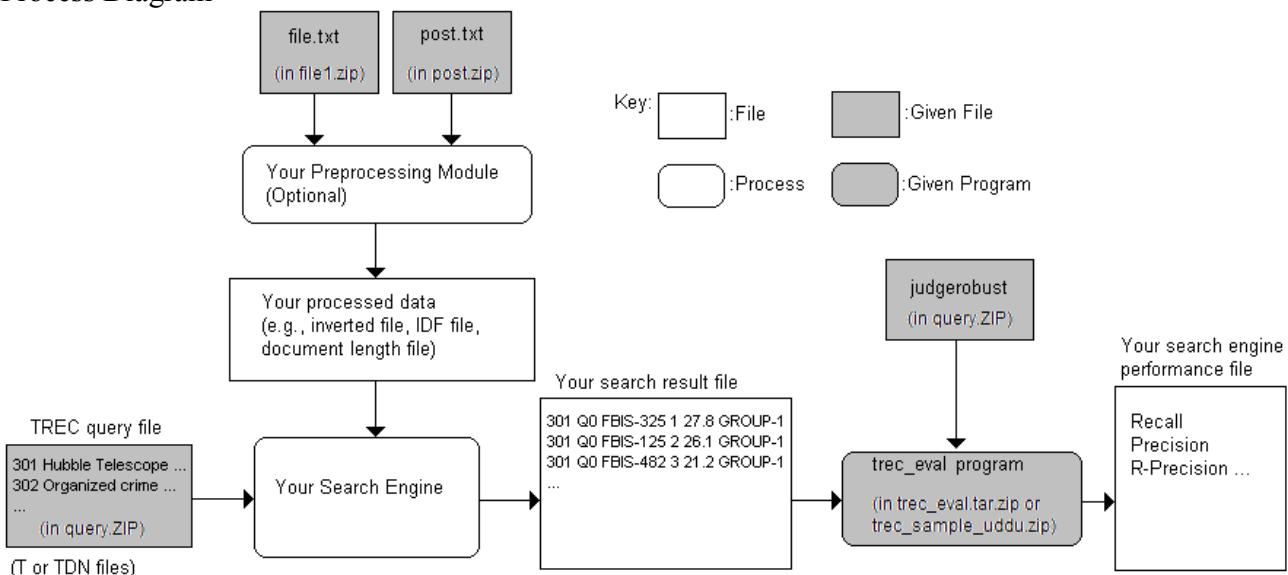
1. Gain hands-on experience in developing a search engine (objective ii and v);
2. Gain hands-on experience in the evaluation of search engines (objective ii and v);
3. Practice communication skill (Attributes of all-roundedness 3 and 4) during demo;

The Problem: You have joined an information retrieval evaluation workshop hosted by me. You need to form a group of size that is between 12 and 16 (inclusive). You should hand in details (i.e., full name, student id, group leader identification and his/her mobile phone contact) to me by **October 17th, 2016**. This workshop seeks to find the most effective and most efficient system that you can design and implement. You are given a set of TREC queries, a set of TREC relevance judgments and some inverted index data. You are required to further develop the system using MS Access or develop another system in C++ (in part or in whole) or another system developed in Java (in part or in whole) to retrieve and rank the given set of documents for the given set of queries. Due to the distribution agreement signed by me with TREC, I cannot distribute the TREC documents to you. However, the inverted index data is available for you to build a search engine.

Given Material (available in Blackboard):

1. prog.zip (C++ programs; example IR system)
2. irproj.zip (programs, readme file, etc.)
3. file.txt (contains the mapping between file identifier and document name used for TREC evaluation)
4. post.txt (contains the inverted index information of the collection)
5. trec_eval.tar.gz (TREC evaluation program and its Windows version)
6. query.zip (Query and relevance judgement files)
 - a. judgerobust: contains the relevance judgment of each query
 - b. queryTDN: long queries
 - c. queryT: short queries

Process Diagram



TREC query format (queryT and queryTDN in query.zip) of each line as follows:

xxxΔ<Query>

where **xxx** are three digits that indicate the query number, **Δ** is the space character and **<Query>** is the content of the query. Figure 1 is an example of a set of queries.

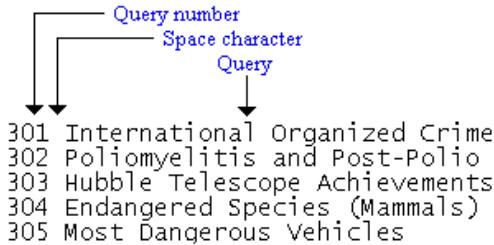


Figure 1: Example TREC queries

TREC Search Result Format (see readme in trec_eval.tar.gz)

The format of the search result (juderobust in query.zip) per line is as follows:

xxxΔQ0Δ<DOC-ID>ΔRΔSΔ<RUN-ID>

where:

XXX are the three digits to specify the query number;
Δ is the space character;
'Q0' is a constant string;
<DOC-ID> is the document identifier;
R is the rank number of the document (1 to 1000 for each query);
S is the similarity score (in descending order);
<RUN-ID> is a string that specifies the run identifier.

Example:

```
301 Q0 FBIS-325 1 27.8 HKPU-1
301 Q0 FBIS-178 2 23.1 HKPU-1
...
...
```

means:

```
query number = 301
document identifier = FBIS-325
rank = 1
similarity score = 27.8
run identifier = HKPU-1
```

Deliverables:

1. Technical Report: 11 points, 1.5 spacing, 20-50 pages. The report should give details of your search engine (e.g. architecture, retrieval models, and ranking techniques), experiments (e.g., experimental design, result and analysis, and conclusion) that you have done, etc. The experiments should report at least the recall-precision curves, the R-precision and the top 10 document precision. When you report the time-efficiency, you should record the speed using the indicated machine in the COMP PC lab. You should also specify the RAM size of the PC and the processor type plus the speed of the system.
2. Search results: Your results should be produced in the TREC search result format (see below) and these results should be stored in the file T.ret and TDN.ret for queryT and queryTDN, respectively. Each query should return at most the top 1000 retrieved documents.
3. CD: It contains your program and relevant data so that I can run your system in one of the COMP PC lab for verification. You should indicate in which COMP PC lab your system is able to run (preferably with an indication of which machine). If I cannot run your program in the indicated COMP PC lab (and machine), marks will be deducted.

Important Dates:

1. Submit your technical report (hardcopy) and CD/DVD (contains your report, powerpoint, search results, etc.) by **November 28th, 2016** to the teaching assistant (Dr Edward Dang, Tel: 2766 7310, Email: cskfdang@comp.polyu.edu.hk) just before the presentation. Late submission will result in mark deduction.
2. You should arrange with the teaching assistant (Contact details are shown in point 1) to demo your system between **November 28th 2016 and December 12th, 2016**. Failure to do so will result in mark deduction. The teaching assistant will run your program using his/her data. If your program does not work, substantial amount of marks will be deducted because we cannot evaluate whether your system is working or not.

Grading

1. In general, if you produce a **basic working system**, then a grade C/C+ will be awarded. You have demonstrated basic understanding of the subject. A basic working system includes the ability to:
 - a. automatically **read in a TREC query file** when a query file is specified by the user and;
 - b. automatically **produce the search results in TREC format** based on some taught retrieval model so that the trec_eval program can be readily used to obtain the performance of the system. If the search result returned 0.0 for precision, recall, or mean average precision, then a grade lower than C might be awarded because it does not show that your system is working.
2. For grade B/B+ in general, you are required to at least achieve the basic working system (i.e., the level of C/C++ grade). In addition, if you are able to implement more than one retrieval model, you will need to **compare and analyze** the merits of one retrieval model over the other. The analysis needs to be substantiated by **measurements** to demonstrate that your analysis is correct. Then, you can be awarded a grade B/B+. Alternatively, if you have only a single retrieval model, you can **analyze the problems** of your retrieval system and **improve** your retrieval system with a new/enhanced method. Then, you can analyze

whether the new/enhance method is performing better than the original method. Again, the analysis needs **measurements** to substantiate the new/enhanced method, and the improvement of the retrieval system needs to be **measured** as well.

3. For grade A/A+ in general, you are required to at least achieve the level of grade B/B+. In addition, you need to have *some innovative elements* in the assignment. The **innovation** needs to be substantial in order for the project to be avoided grade A/A+. Hence, trivial innovations do not count. Example of substantial innovation includes inventing a new retrieval model that can substantially improve the retrieval effectiveness of your basic system or improved system. Again, you need to substantiate that your innovation is effective by gathering **evidence**.
4. **You can use the programs provided to you by this subject for this assignment, but you are not allowed to use any open source IR programs.** You cannot share programs with other groups unless these programs are provided by us. Fail grades will be awarded to those who are found to participate in plagiarism or in copying other group's programs.
5. For project management, you need to **put down the percentage of contribution in the report for each group member**. If the group does not put down the percentage of contribution in the report, the grade for all members of the group will be lowered. If you have **0% contribution**, then your grade may be significantly lower than the grade of your group mates. The accepted report is based on the one handed to me just before the presentation (by the group leader). If you hand in it late, you then let me to decide which report I want to take as the final one. If you have arguments about the contribution, then whoever has 0% contribution should write down his/her contribution to the assignment, and together with the group leader submit the contribution to the Teaching Assistant (TA) within one week after the presentation. If the group leader is not present, the person with 0% contribution should provide evidence to the TA that the group leader is being contacted but did not reply (e.g., by email with sufficient time to reply, e.g. within 2 days). If the group leader did not respond and did not come to the TA with the concerned person, then the grade of the group leader may be lowered. The responsibility of the group leader is to verify in the presence of the TA that the claimed contribution by the concerned person is true. The person with 0% contribution should **provide evidence** to support that (s)he made the contribution (e.g., email that contains **part of the report, minutes of meeting** indicating whose idea was used, etc). If the dispute about the contribution is not settled, then the grade of the entire group may be lowered because this indicates that the project management is not done well.
6. **To show that you have done some project management, you should include the project plan, schedule, requirements identification, task identification, job allocation and distribution (who is assigned with what tasks), etc.** in the report.

Initial Set Up Procedure:

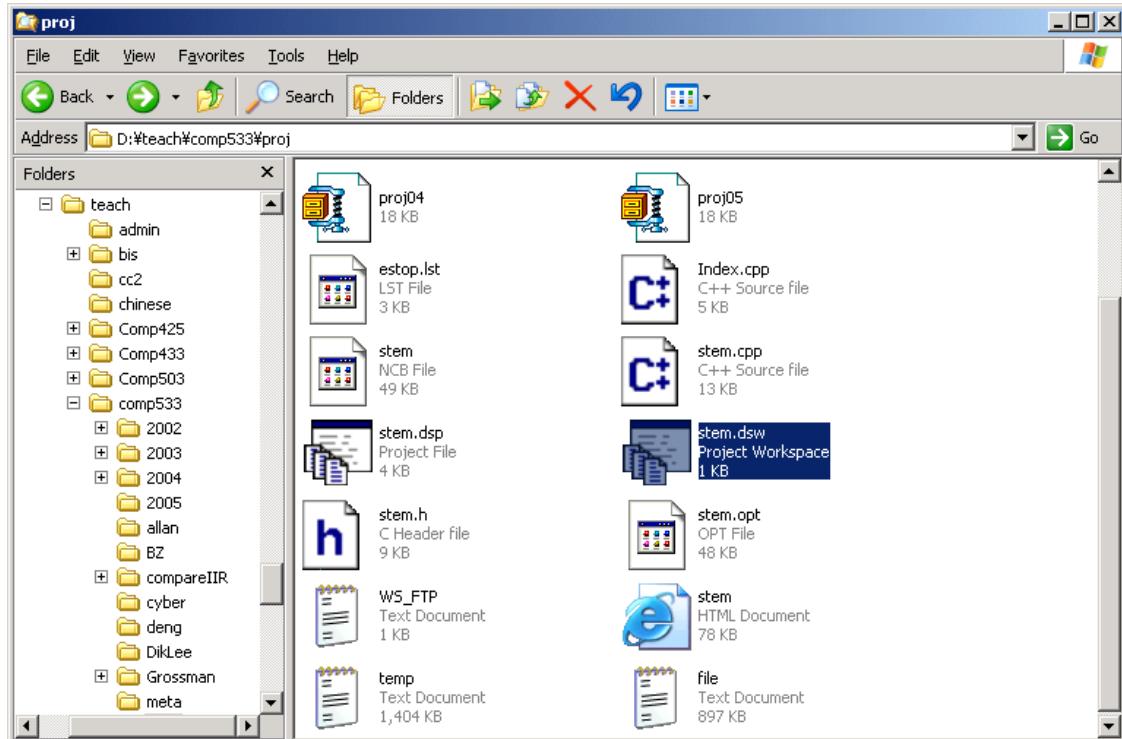
The first 10 steps are for reference only, because we can not release our TREC documents to you. The processing is already done for you, but it shows how the data are obtained so that you have a complete picture of the process.

From step 11, you can follow it and use MS Access to build a search engine (although not recommended because it is too slow) or use it to preprocess data.

Step 1. Download the IRProject file (For illustration) [done]

Step 2. Unzip the IRProject file (For illustration) [done]

Step 3. In MS Windows with Visual C++ installed, click on **stem.dsw** (which is in reverse video in the following figure) (For illustration) [done]



Step 4: You should see MS Visual C++ development environment as follows: (For illustration) [done]

```
stem - Microsoft Visual C++ - [Index.cpp]
File Edit View Insert Project Build Tools Window Help
(Globals) (All global members) main
Workspace 'stem': 1 project(s)
stem files
Index.cpp
stem.cpp
stem.h
// Information Retrieval Assignment //
// Robert Luk @ 2005 //
// All rights reserved //
// Department of Computing //
// The Hong Kong Polytechnic University //
// //
// This program should only be used for //
// the information retrieval subjects //
// (COMP425, COMP433 and COMP5324) of the //
// Department of Computing, The Hong Kong //
// Polytechnic University. //

```

Step 5: Scroll down the Index.cpp file until you see the dir variable declared as follows: (For illustration) [done]

```
stem - Microsoft Visual C++ - [Index.cpp *]
File Edit View Insert Project Build Tools Window Help
(Globals) (All global members) main
Workspace 'stem': 1 project(s)
stem files
Index.cpp
stem.cpp
stem.h
int i;
// hard code your directory that stores the data files
char * dir = "dir/b/s YYteachYYcomp533YYdataYY > temp.txt\n";
system(dir);
fp = fopen("temp.txt", "rb");
ffp = fopen("file.txt", "w");
ofp = fopen("post.txt", "w");
LoadStop(".YYstop.lst");
// Start processing each file in the temp.txt file
while (fgets(line, 1000, fp) != NULL) [
    // Get the file name.
```

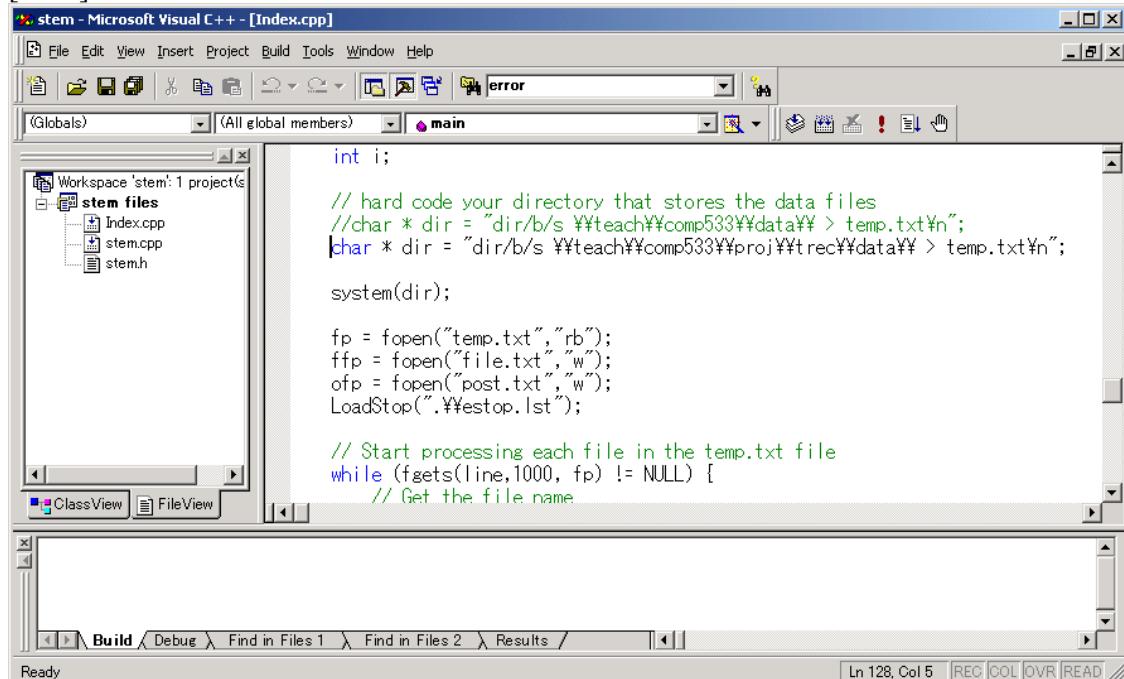
Step 6: Change the directory specification of the dir variable to the directory that stores your data. (For illustration) For example, we have changed the directory to:

```
\\"teach\\comp533\\proj\\trec\\data\\
```

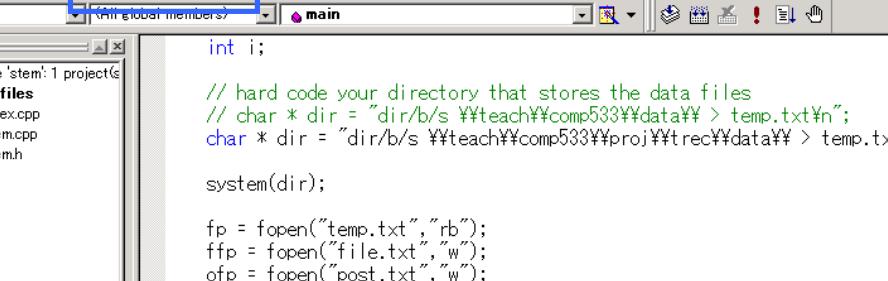
and we have specify the value of the dir variable as follows:

```
char * dir = "dir/b/s \\"teach\\comp533\\proj\\trec\\data\\\" > temp.txt\\n";
```

[done]



Step 7: Compile the program. Click the Build menu (top) and click the Rebuild All option in the pull down menu. (For illustration) [done]



The screenshot shows the Microsoft Visual C++ IDE interface. The title bar reads "stem - Microsoft Visual C++ - [Index.cpp]". The menu bar includes File, Edit, View, Insert, Project, Build, Tools, Window, and Help. A blue box highlights the "Project" menu item. The toolbar contains various icons for file operations like Open, Save, and Build. The status bar at the bottom shows "Ln 127, Col 8" and "REC COL OVR READ".

int i;

```
// hard code your directory that stores the data files
// char * dir = "dir/b/s %%teach%%comp533%%data%% > temp.txt\n";
char * dir = "dir/b/s %%teach%%comp533%%proj%%rec%%data%% > temp.txt\n";

system(dir);

fp = fopen("temp.txt", "rb");
fpr = fopen("file.txt", "w");
ofp = fopen("post.txt", "w");
LoadStop(".%%estop.lst");

// Start processing each file in the temp.txt file
while (fgets(line, 1000, fp) != NULL) {
    // Get the file name.
```

The screenshot shows the Microsoft Visual C++ 6.0 IDE interface. The title bar reads "stem - Microsoft Visual C++ - [Index.cpp]". The menu bar includes File, Edit, View, Insert, Project, Build, Tools, Window, Help. The toolbar has icons for New, Open, Save, Print, Find, Replace, Cut, Copy, Paste, Undo, Redo, and others. The status bar at the bottom shows "Ln 386, Col 1 REC COL OVR READ".

Workspace: stem: 1 project(s)

- stem files
 - Index.cpp
 - stem.cpp
 - stem.h

Code Editor:

```
int i;

// hard code your directory that stores the data files
// char * dir = "dir/b/s YYteachYYcomp533YYdataYY > temp.txt\n";
char * dir = "dir/b/s YYteachYYcomp533YYprojYYtrecYYdataYY > temp.txt\n";

system(dir);

fp = fopen("temp.txt", "rb");
ffp = fopen("file.txt", "w");
ofp = fopen("post.txt", "w");
LoadStop("./estop.lst");

// Start processing each file in the temp.txt file
while (fread(lin, 1000, 1, fp) != NULL) {
```

Status Bar: warning C4786: 'std::_Tree<std::basic_string<char, std::char_traits<char>, std::allocator<char> >, std::pair<char> >, int, std::less<std::basic_string<char, std::char_traits<char>, std::allocator<char> > >, std::allocator<char> >' was truncated to '255' characters in the debug information
Linking...

stem.exe - 0 error(s), 101 warning(s)

Step 8: After compiling the program, copy the compiled program **stem.exe** (which is stored in the child directory **Debug** of the current directory) to the current directory. (For illustration) [done]

```

Select Command Prompt
Volume in drive D is DATA
Volume Serial Number is 1950-12EE

Directory of D:\teach\comp533\proj\Debug

26/01/2005  02:51 PM    <DIR>      .
26/01/2005  02:51 PM    <DIR>      ..
26/01/2005  04:53 PM    <DIR>      test
27/01/2005  09:51 AM           91,136 vc60.idb
27/01/2005  09:51 AM           118,784 vc60.pdb
27/01/2005  09:51 AM           198,864 Index.obj
27/01/2005  09:51 AM           2,596,036 stem.pch
27/01/2005  09:51 AM           51,045 stem.obj
27/01/2005  09:51 AM           321,332 stem.llk
27/01/2005  09:51 AM           225,333 stem.exe
27/01/2005  09:51 AM           541,696 stem.pdb
               8 File(s)   4,144,226 bytes
               3 Dir(s)  9,314,893,824 bytes free

D:\teach\comp533\proj\Debug>copy stem.exe ...

```

After copying:

```

Select Command Prompt
27/01/2005  09:44 AM           48,640 stem.opt
27/01/2005  09:54 AM           322,560 readme.doc
27/01/2005  09:51 AM           225,333 stem.exe
               16 File(s)  361,705,733 bytes
               5 Dir(s)  9,314,893,824 bytes free

D:\teach\comp533\proj>dir /w
Volume in drive D is DATA
Volume Serial Number is 1950-12EE

Directory of D:\teach\comp533\proj

[.]      [...]      stem.cpp      stem.h      [Debug]      stem.ncb
stem.plg  stem.dsp  stem.dsw  post.txt  estop.lst  [data]
[trec]  proj04.ZIP  proj05.zip  WS_FTP.LOG  file.txt  Index.cpp
stem.opt  readme.doc  stem.exe
               16 File(s)  361,705,733 bytes
               5 Dir(s)  9,314,893,824 bytes free

D:\teach\comp533\proj>stem.exe

```

Step 9: After making sure that the current directory has the **estop.lst** file (i.e., the stop word list), you can run the program **stem.exe** as above. When **stem.exe** is running, it will report it has processed every 100 documents as in the screen shot below. (For illustration) [done]

```

Command Prompt - stem.exe
D:\teach\comp533\proj>stem.exe
[ 0]Reading file:D:\teach\comp533\proj\trec\data\FBIS3-10291
[ 100]Reading file:D:\teach\comp533\proj\trec\data\FBIS3-38725
[ 200]Reading file:D:\teach\comp533\proj\trec\data\FBIS3-58778
[ 300]Reading file:D:\teach\comp533\proj\trec\data\FBIS4-24424
[ 400]Reading file:D:\teach\comp533\proj\trec\data\FBIS4-43862
[ 500]Reading file:D:\teach\comp533\proj\trec\data\FBIS4-65419
[ 600]Reading file:D:\teach\comp533\proj\trec\data\FR940513-2-00009
[ 700]Reading file:D:\teach\comp533\proj\trec\data\LA082590-0058
[ 800]Reading file:D:\teach\comp533\proj\trec\data\FT924-9772
[ 900]Reading file:D:\teach\comp533\proj\trec\data\FT911-5038
[1000]Reading file:D:\teach\comp533\proj\trec\data\FBIS3-42477
[1100]Reading file:D:\teach\comp533\proj\trec\data\FBIS3-9183
[1200]Reading file:D:\teach\comp533\proj\trec\data\FBIS4-54491
[1300]Reading file:D:\teach\comp533\proj\trec\data\FBIS4-26677
[1400]Reading file:D:\teach\comp533\proj\trec\data\LA111489-0151
[1500]Reading file:D:\teach\comp533\proj\trec\data\FT922-4763
[1600]Reading file:D:\teach\comp533\proj\trec\data\FBIS4-7732
[1700]Reading file:D:\teach\comp533\proj\trec\data\FT942-12967

```

Step 10: The program, stem.exe, produces two files: (For illustration)

post.txt: this file contains the posting data which has the following format

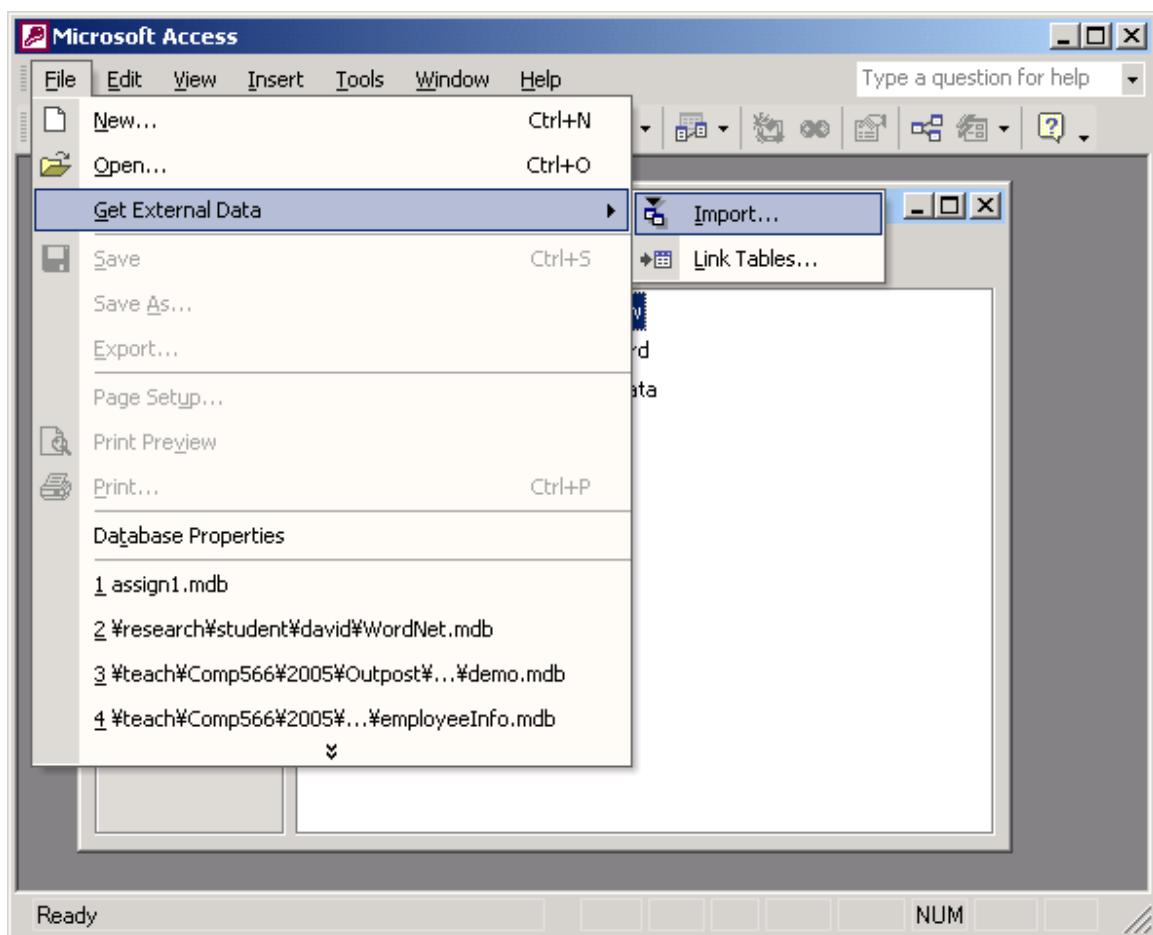
field	name	term	file-id	logical word position
example		chinesE	2	13

file.txt: this file contains the file id, name and path data. The format is as follows:

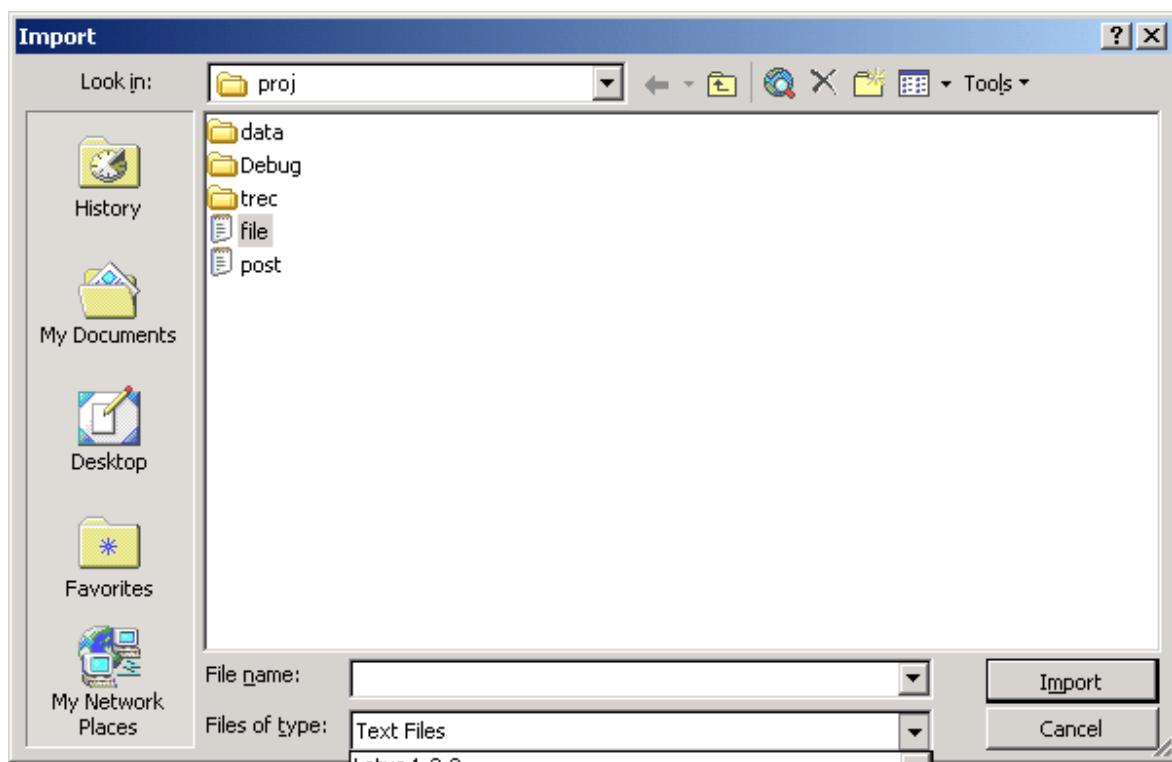
field	name	fileid	doclen	status	docid	path
example	13	377	@	FBIS3-14196	D:\teach\comp533\proj\trec\data\FBIS3-14196	

[done]

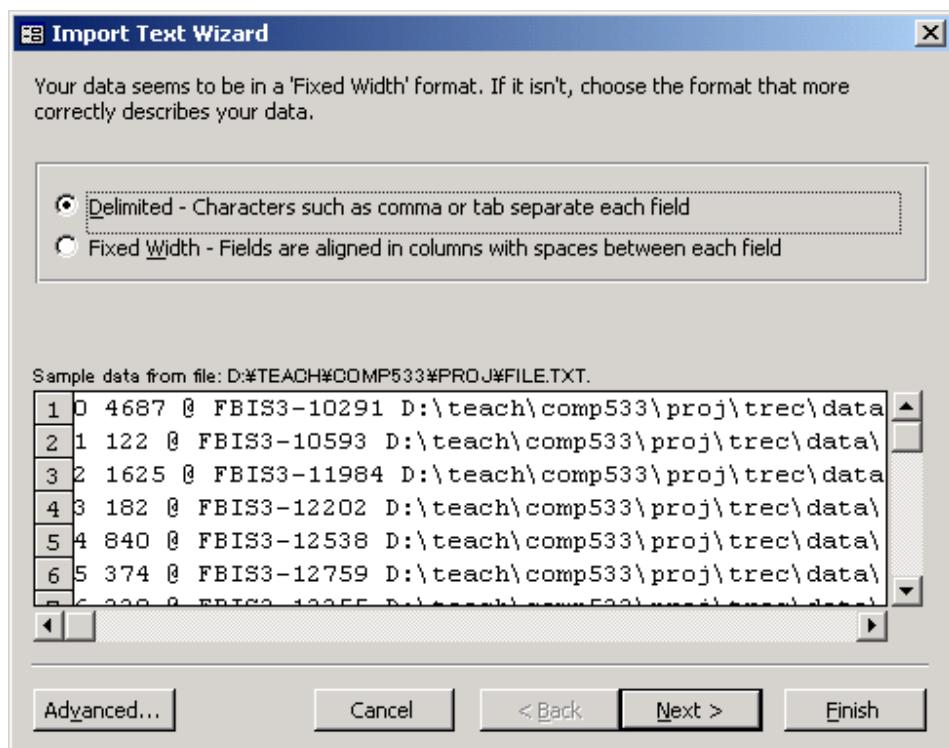
Step 11: Import these two files into your new MS Access database into two tables: post table for data in post.txt and file table for data in file.txt



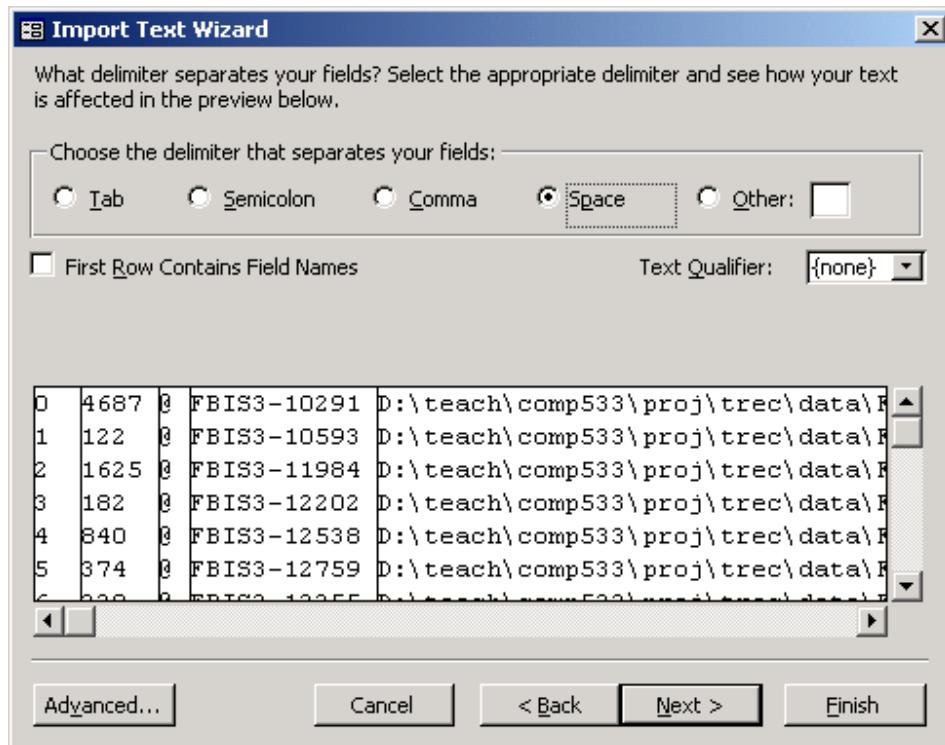
Import based on text files:



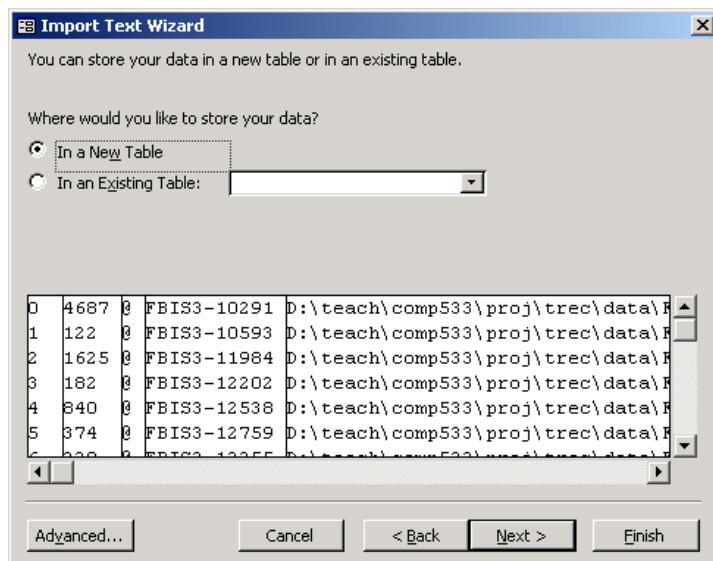
Step 12: Click on **Delimited** and then **Next**.



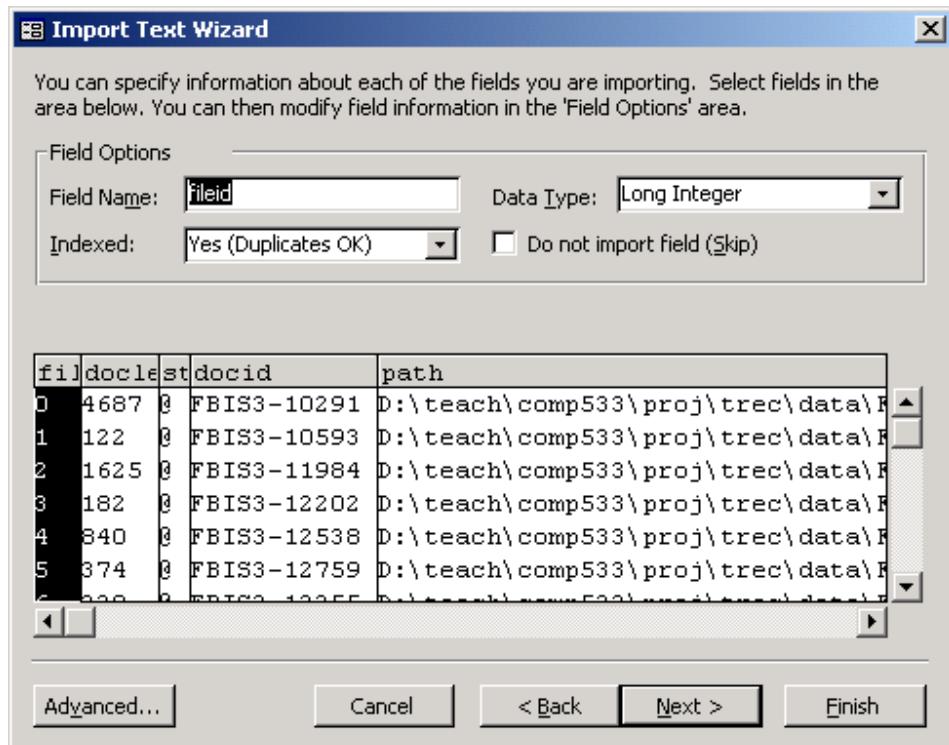
Step 13: Click the Space radial button and press Next.



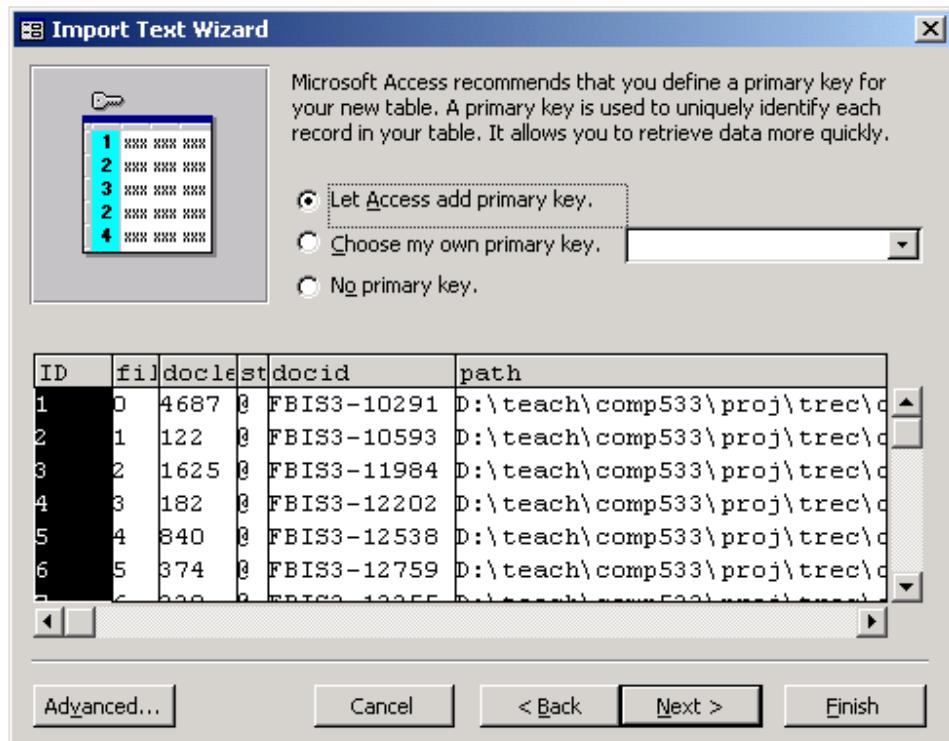
Step 14: Put the data into the new table (file.txt) and press Next.



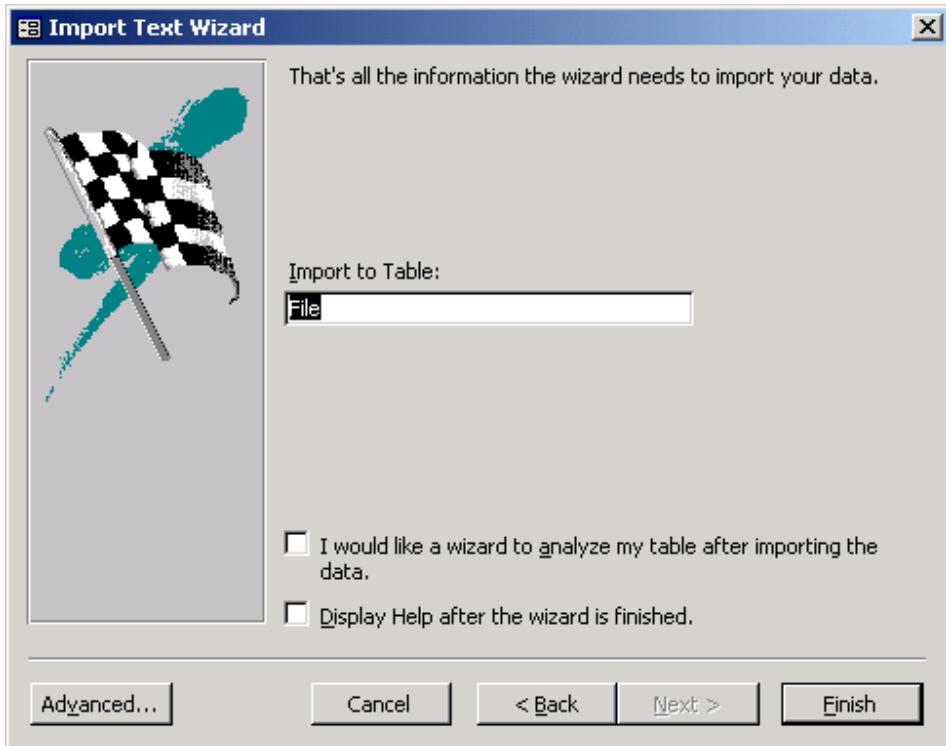
Step 15: Enter the names of each file.



Step 16: Set the primary key to **fileid**.



Step 17: Enter the name of the able as **File**.



Step 19: After importing the post.txt and file.txt, your database should look like the following:

The Microsoft Access interface shows two tables: "Post" and "File".

Post : Table

ID	term	fileid	wpos
1	doC	0	0
2	docnO	0	1
3	fbI	0	2
4	docnO	0	3
5	hT	0	4
6	drlaT	0	5
7	hT	0	6

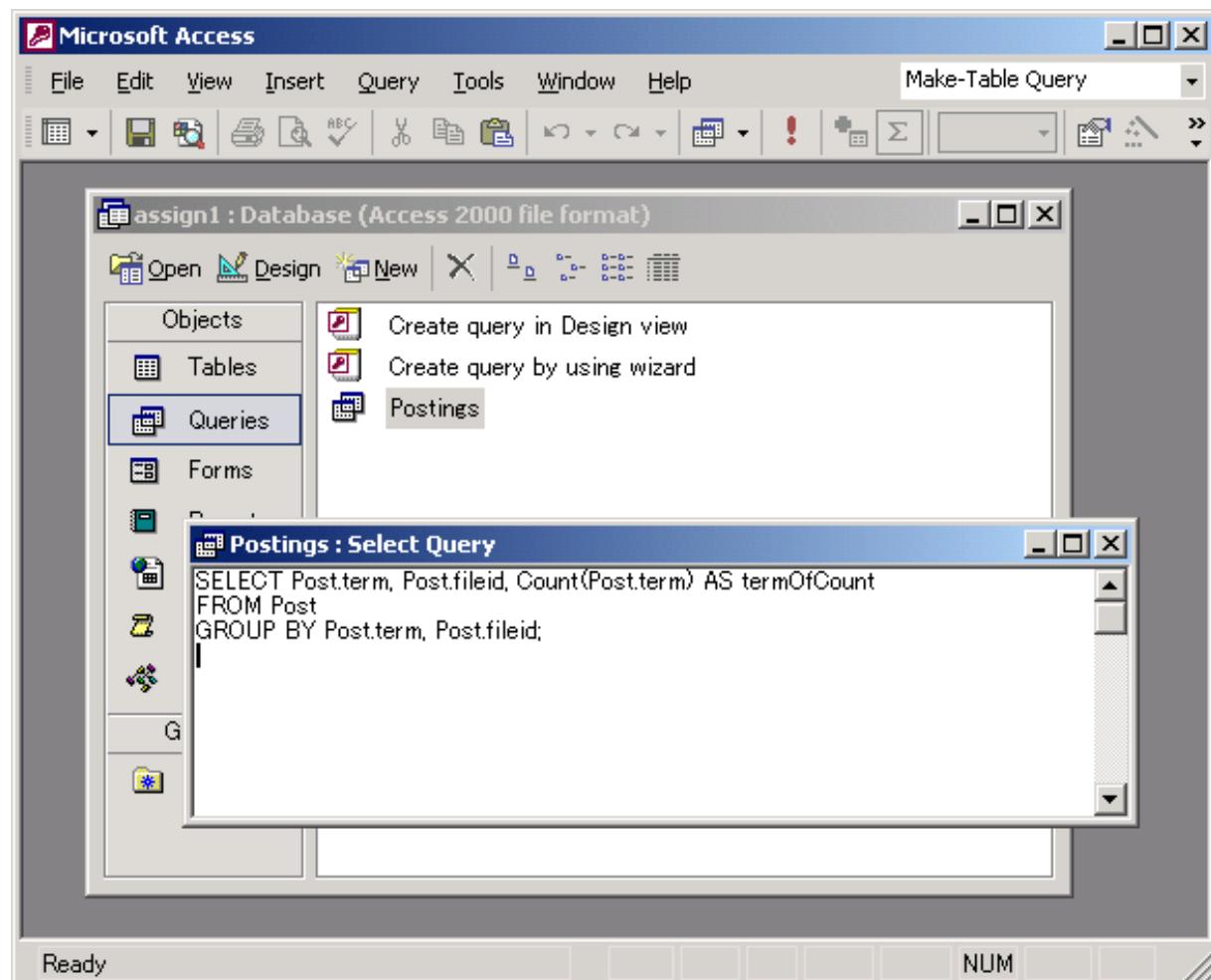
File : Table

fileid	doclen	status	docid	path
0	0 @		FBIS3-10291	D:\teach\comp533\proj\trec\data\FBIS3-10291
1	0 @		FBIS3-10593	D:\teach\comp533\proj\trec\data\FBIS3-10593
2	0 @		FBIS3-11984	D:\teach\comp533\proj\trec\data\FBIS3-11984
3	0 @		FBIS3-12202	D:\teach\comp533\proj\trec\data\FBIS3-12202
4	0 @		FBIS3-12538	D:\teach\comp533\proj\trec\data\FBIS3-12538
5	0 @		FBIS3-12759	D:\teach\comp533\proj\trec\data\FBIS3-12759
6	0 @		FBIS3-13355	D:\teach\comp533\proj\trec\data\FBIS3-13355
7	0 @		FBIS3-16004	D:\teach\comp533\proj\trec\data\FBIS3-16004

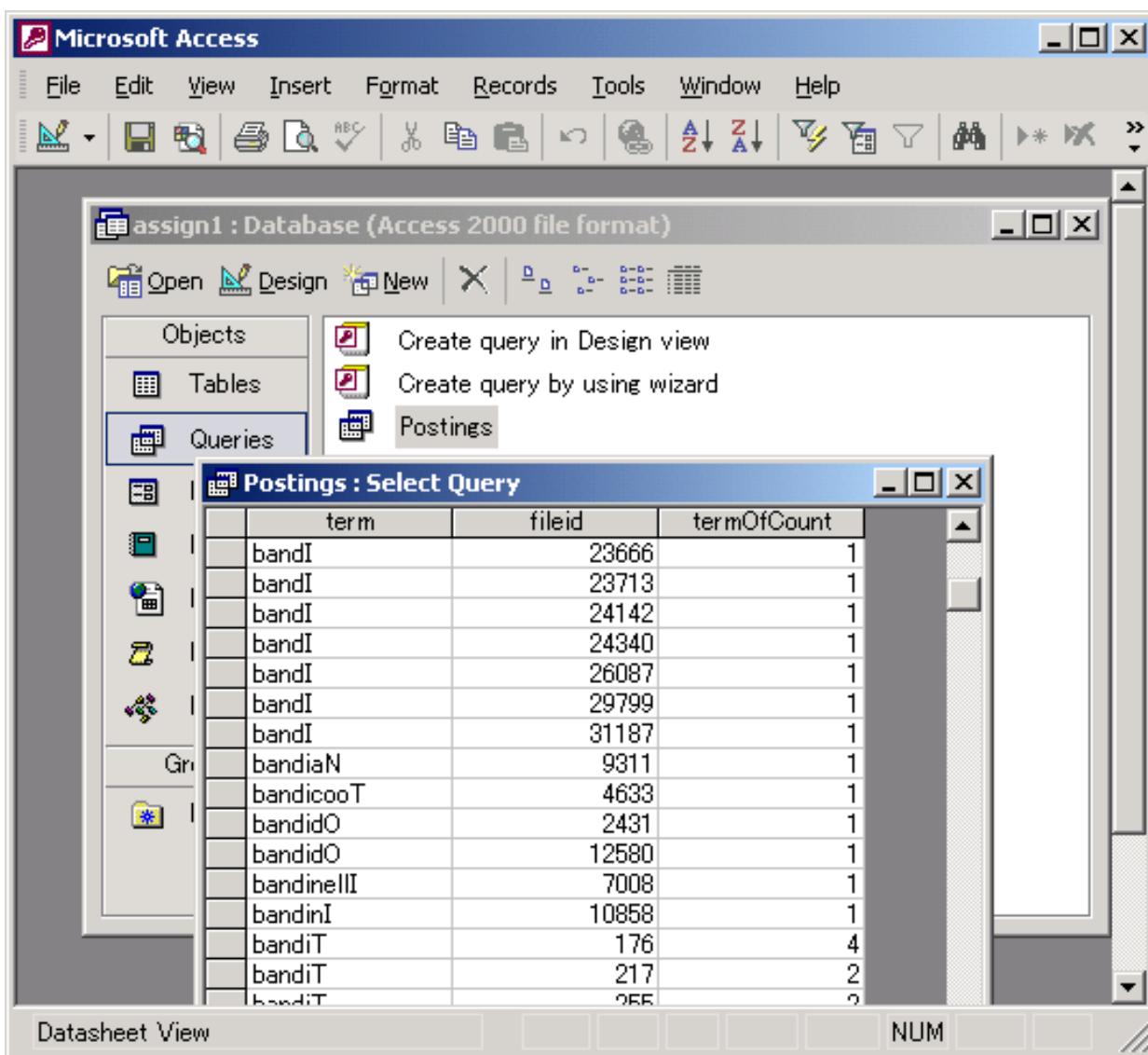
The file table has 64814 records.

Step 20: Create a query (called Postings) to count the within-document term frequency. First, you need to create a query as follows in SQL:

```
SELECT Post.term, Post.fileid, Count(Post.term) AS termOfCount
FROM Post
GROUP BY Post.term, Post.fileid;
```



Step 21: Compile the query:



Step 22: Based on the data (e.g. file.txt + post.txt or the MS Access file that stores these data) you obtained, you should develop an IR system using this information. There are several possible ways that you can develop the system as follows:

- (a) use MS Access and possibly embedding some VBA;
- (b) use MS Access as the database server and write a client to transfer the posting information to the client during retrieval. Your client program may be written in JAVA or C++ using some database connection utility like ODBC;
- (c) develop a customized JAVA or C++ program that directly reads data from files.txt and post.txt (or some derived data). In this case, don't use any database package to help
- (d) use another method but you need to make sure that the computing lab has the same package (e.g. using another database software) so that you can do the demonstration.

You are not allowed to use any open source search engine.

Step 23: When you input a query to the search engine, you need to stem the words in the query the same was as the index terms are stemmed. The stemmer program is already available in C++. This program is stored in stem.cpp and stem.h, based on the Porter stemming algorithm. To use this program in C++,

- (a) declare that you are using this stemmer in the program by including the header file and declaring a stemmer object:

```
#include <stem.h>
stemmer StemProc;
```

- (b) call the member function Stem(.) to stem the word;

```
char word[1000];
strcpy(word, "organs")
StemProc.Stem(word); // word contains "orgaN"
```

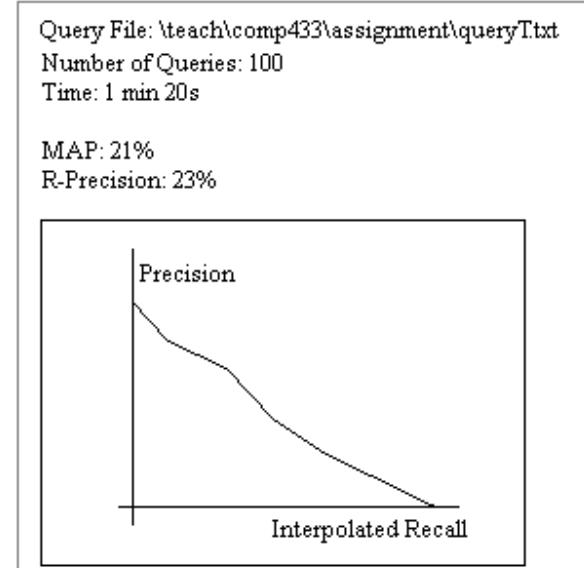
If you develop your JAVA program, you need to download a JAVA implementation of the Porter stemming algorithm from the web. In this case, you need to modify the JAVA program so that after stemming, the last character of the stemmed word is a capital letter (e.g., the last letter is in upper case for “orgaN”). This is because our stemmer capitalized the last letter of the index terms. Otherwise, it will be difficult for you to match the index terms with the query terms.

Step 24: A possible interface for the system may be:

A possible system **input interface**:

Query File:	<input type="text"/>	<input type="button" value="Browse"/>
Retrieval model	Term Weight Normalization	
<input type="checkbox"/> Boolean	<input type="checkbox"/> Yes	
<input type="checkbox"/> Vector Space	<input type="checkbox"/> No	
<input type="checkbox"/> Fuzzy Boolean		

A possible system **output interface**



NB: the TREC evaluation program (trec_eval.tar) produces the MAP, R-precision, interpolated recall-precision curve based in the ranked output of your search engine and the file (i.e., judgerrobust) containing the correct answers.