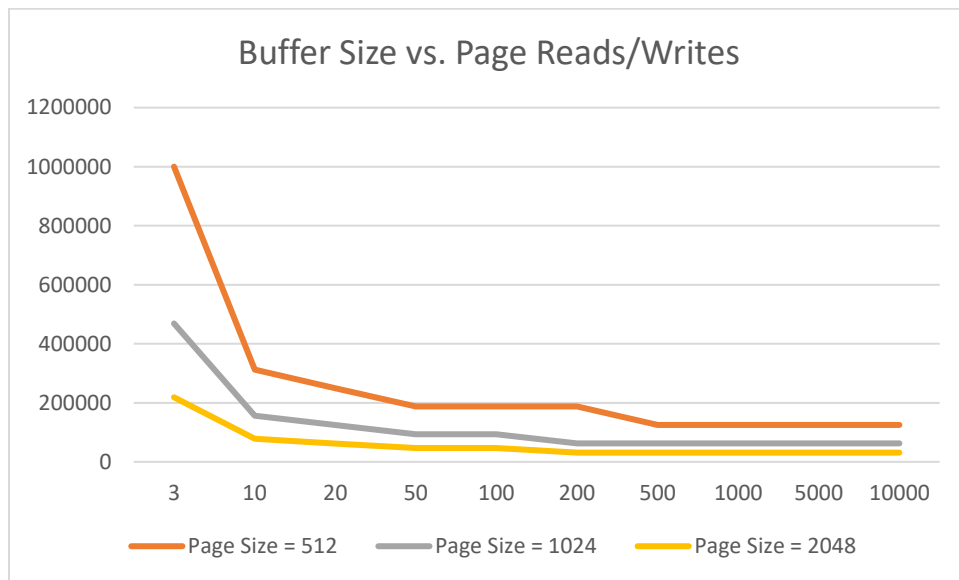


Part one Answers:

Total page reads/writes										
Page Size\B	3	10	20	50	100	200	500	1000	5000	10000
512	1000000	312500	250000	187500	187500	187500	125000	125000	125000	125000
1024	468750	156250	125000	93750	93750	62500	62500	62500	62500	62500
2048	218750	78125	62500	46875	46875	31250	31250	31250	31250	31250

Total Passes (unit in pages)										
Page size\B	3	10	20	50	100	200	500	1000	5000	10000
512	16	5	4	3	3	3	2	2	2	2
1024	15	5	4	3	3	2	2	2	2	2
2048	14	5	4	3	3	2	2	2	2	2



Observations: From the line chart, we see that the number of pages read (y-axis) diminish as the number of buffer pages (x-axis) increases. These empirical results match the equation provided in class, which is $\log_{B-1} \frac{N}{B}$, as B increases, the number gets smaller.

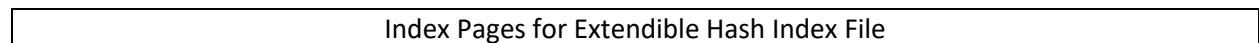
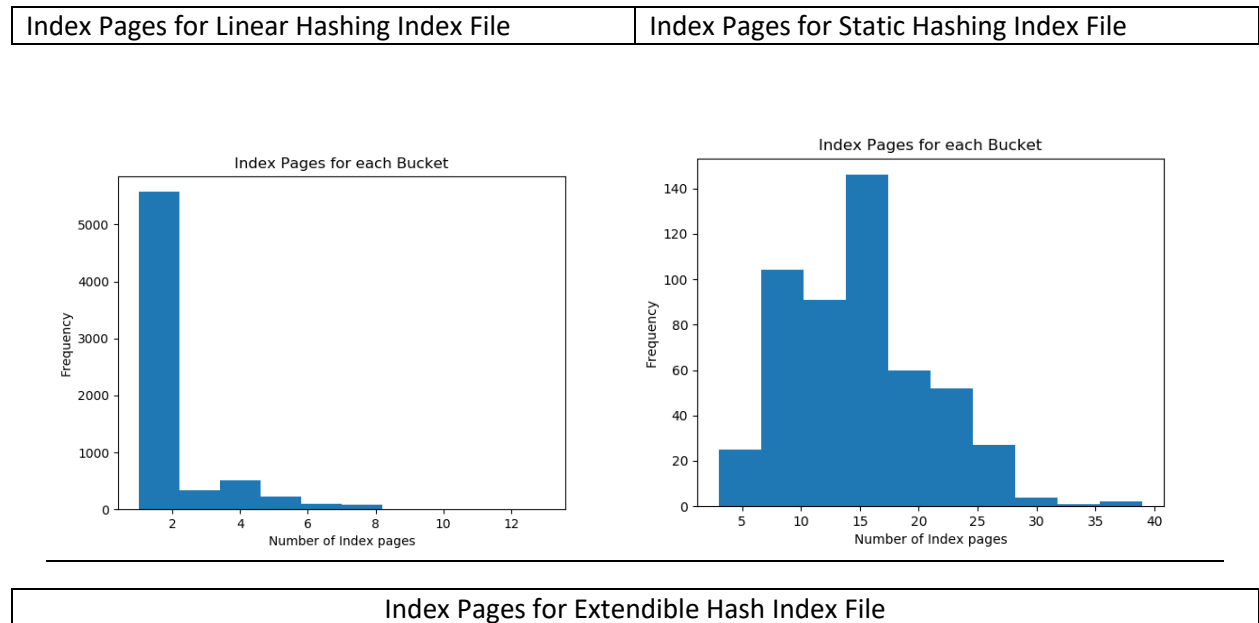
Name	Date modified	Type	Size
output_size_512_B_3	3/6/2019 11:20 AM	Data Base File	31,250 KB
output_size_512_B_10	3/6/2019 11:25 AM	Data Base File	31,250 KB
output_size_512_B_20	3/6/2019 11:27 AM	Data Base File	31,250 KB
output_size_512_B_50	3/6/2019 11:29 AM	Data Base File	31,250 KB
output_size_512_B_100	3/6/2019 11:31 AM	Data Base File	31,250 KB
output_size_512_B_200	3/6/2019 11:33 AM	Data Base File	31,250 KB
output_size_512_B_500	3/6/2019 11:34 AM	Data Base File	31,250 KB
output_size_512_B_1000	3/6/2019 11:35 AM	Data Base File	31,250 KB
output_size_512_B_5000	3/6/2019 11:37 AM	Data Base File	31,250 KB
output_size_512_B_10000	3/6/2019 11:39 AM	Data Base File	31,250 KB
output_size_1024_B_3	3/6/2019 11:47 AM	Data Base File	31,250 KB
output_size_1024_B_10	3/6/2019 11:50 AM	Data Base File	31,250 KB
output_size_1024_B_20	3/6/2019 11:52 AM	Data Base File	31,250 KB
output_size_1024_B_50	3/6/2019 11:54 AM	Data Base File	31,250 KB
output_size_1024_B_100	3/6/2019 11:56 AM	Data Base File	31,250 KB
output_size_1024_B_200	3/6/2019 11:57 AM	Data Base File	31,250 KB
output_size_1024_B_500	3/6/2019 11:58 AM	Data Base File	31,250 KB
output_size_1024_B_1000	3/6/2019 11:59 AM	Data Base File	31,250 KB
output_size_1024_B_5000	3/6/2019 12:00 PM	Data Base File	31,250 KB
output_size_1024_B_10000	3/6/2019 12:01 PM	Data Base File	31,250 KB
output_size_2048_B_3	3/6/2019 12:09 PM	Data Base File	31,250 KB
output_size_2048_B_10	3/6/2019 12:13 PM	Data Base File	31,250 KB
output_size_2048_B_20	3/6/2019 12:15 PM	Data Base File	31,250 KB
output_size_2048_B_50	3/6/2019 12:16 PM	Data Base File	31,250 KB
output_size_2048_B_100	3/6/2019 12:18 PM	Data Base File	31,250 KB
output_size_2048_B_200	3/6/2019 12:19 PM	Data Base File	31,250 KB
output_size_2048_B_500	3/6/2019 12:19 PM	Data Base File	31,250 KB
output_size_2048_B_1000	3/6/2019 12:20 PM	Data Base File	31,250 KB
output_size_2048_B_5000	3/6/2019 12:21 PM	Data Base File	31,250 KB
output_size_2048_B_10000	3/6/2019 12:23 PM	Data Base File	31,250 KB

```

467
468 if __name__ == "__main__":
469     default = "/names.db"
470     inputFile = "names.db"
471     outputFile = "output"
472     dbFile = "output"
473     pageSize = [1512, 1024, 2048]
474     Bn = [1, 10, 100, 1000, 10000, 100000, 1000000, 10000000]
475     # redirects the outputs
476     # open the
477     outputFile = open("./results/externalSorting_outputs.txt", "w")
478     sys.stdout = outputFile
479     for pageSize in pageSize:
480         for b in Bn:
481             # if the file names for different combinations
482             outputName = outputFile + ".size {}".format(pageSize) + ".b {}".format(b) + ".db"
483             # BDB files have each combination of sorting too a clean startup
484             countStart(inputFile, outputName)
485             inputFile = open("./{}".format(inputFile), "rb")
486             outputFile = open("./results/externalSorting_SortedFiles/{}".format(outputName), "wb")
487             externalSorter = ExternalSorter(inputFile, RECORDSIZE, Field, inputFile, outputSize, dbFile, inputFile, outputName)
488             inputFile.close()
489             outputFile.close()
490             # resume the file if the inputFile name is the output file
491             if inputFile.name != dbFile:
492                 trueOutputName = inputFile.name
493                 trueOutputName = outputFile.name
494                 temp = "temp"
495                 os.rename(inputFile.name, temp)
496                 os.rename(outputName, trueOutputName)
497                 os.rename(temp, trueOutputName)
498             # handle statistics
499             statTracker = externalSorter.getStatTracker()
500             # output the statistics to the file
501             print("Page Size: {} - Buffer Size: {}".format(pageSize, b))
502             print(statTracker)
503             print("")
504             outputFile.close()
505
506

```

Part two answers:

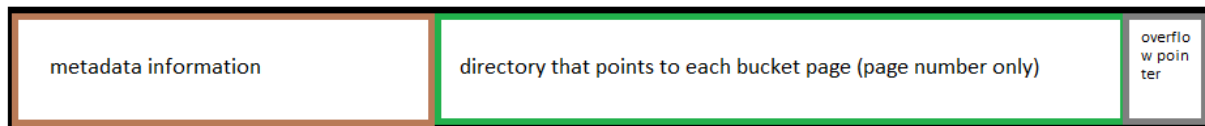


Outputs from the program:

	number of buckets	Regular index pages	Overflow pages
Static hash index	512	512	7084
Extendible hash index	8814	9843	4393
Linear hash index	6845	6845	4967

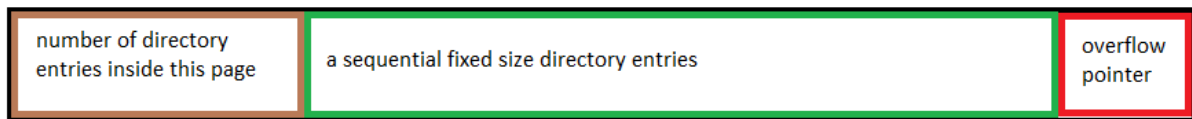
Static Hash Index File **header page** Format contains the following information in order:

- Page size – 2 bytes (determine each page size of the index file)
- Index file type – 1 byte (to identify the index type of this index file)
- Fixed number of buckets – 1 byte (number of fixed buckets in this index file)
- Overflow page pointer size – 1 byte (number of bytes to read for the overflow page number)
- Record entry size – 1 byte (the size of each record entry in index and overflow pages)
- Bytes reserved to store number of record entries are stored in index page (exclude directory page) – 1 byte
- Directory entry size – 2 bytes (the size of each directory entry)
- A directory of buckets that points to the index page for the rest of the header page
- Overflow pointer for the next directory page – 2 bytes (points to the next directory page)



Static Hash index File **directory page** Format contains the following information in order:

- Number of directory entries in this page – 2 bytes (number of fixed size directory entries in this page)
- A list of directory entries (a sequential list of fixed size directory entries)
- Overflow directory pointer to the next directory page – 2 bytes (the page number for the next directory page)



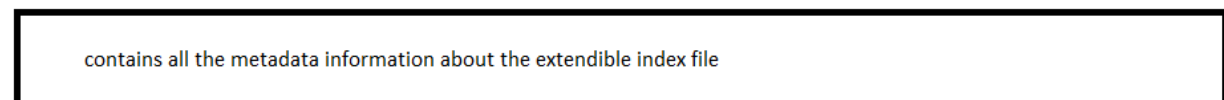
Static Hash index File **index and overflow page** format contains the following information in order:

- Number of records in this page – 2 bytes (number of fixed size record entries in this page)
- All the record entries (a sequential list of fixed size record entries)
- At the end, there is an overflow page pointer – 2 bytes (the page number for the overflow page of this bucket)



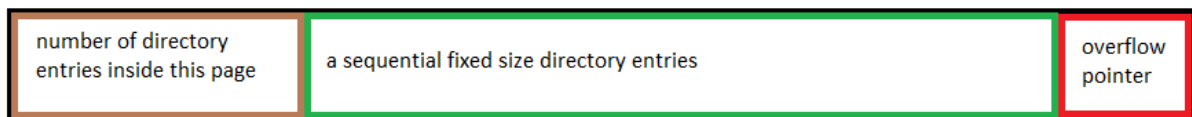
Extendible Hash Index File **header page** Format contains the following information in order:

- Page size – 2 bytes (determine each page size of the index file)
- Index file type – 1 byte (to identify the index type of this index file)
- First directory page start location – 1 byte (the page number of the first directory page)
- Global depth – 1 byte (the global depth is needed for the hash function)
- Overflow page pointer size – 1 byte (number of bytes to read to interpret the overflow page number)
- Directory entry size – 1 byte (size of each directory entry in directory pages)
- Record entry size – 1 byte (size of each record entry in index and overflow pages)
- Bytes reserved to store the local depth for each bucket page – 1 byte (number of bytes to read to interpret the local depth of the index page)
- Bytes reserved to store number of record entries are stored in index page (exclude directory page) – 1 byte (number of bytes to read to know how many records in this page)



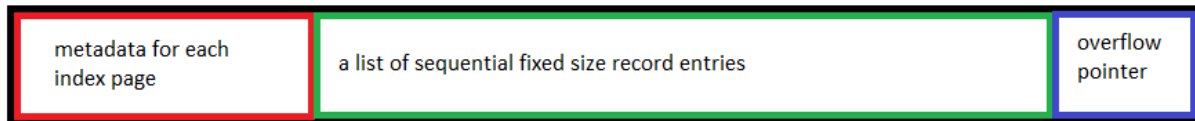
Extendible Hash Index File **directory page** Format contains the following information in order:

- Number of directory entries in this page – 2 bytes (number of fixed size directory entries in this page)
- A list of directory entries (a sequential list of fixed size directory entries)
- Overflow directory pointer to the next directory page – 2 bytes (the page number for the next directory page)



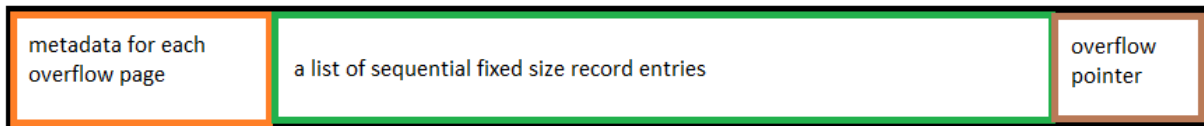
Extendible Hash Index File **index page of each bucket** contains the following information in order:

- Local depth of this bucket – 2 bytes (the local depth of this index page, this is required to determine whether to double the directories)
- Number of record entries in this page – 2 bytes (the number of record entries currently in this page)
- A sequential fixed size list of record entries (a list of fixed size record entries)
- Overflow page pointer – 2 bytes (the overflow page number of the next overflow page of this bucket)



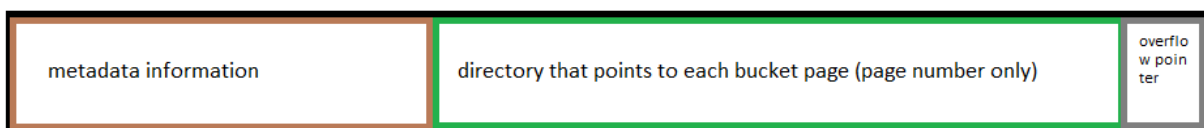
Extendible Hash Index File **overflow page of each bucket** contains the following information in order:

- Number of records in this page – 2 bytes (number of fixed size record entries in this page)
- All the record entries (a sequential list of fixed size record entries)
- At the end, there is an overflow page pointer – 2 bytes (the page number for the overflow page of this bucket)



Linear Hash Index File **header page** Format:

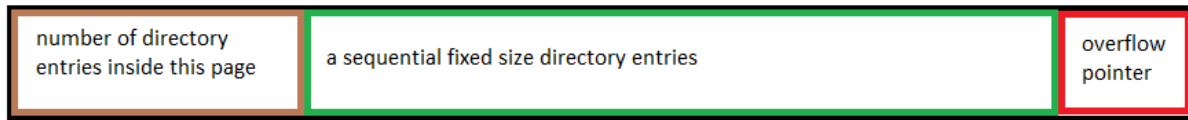
- Page size – 2 bytes (determine each page size of the index file)
- Index file type – 1 byte (to identify the index type of this index file)
- Initial number of buckets – 1 byte (the initial number of buckets to begin with, need for the hash function)
- Current level – 1 byte (the current level, it is need for the hash function)
- Next bucket pointer – 1 byte (the bucket that the *next* pointer points to, needed to determine how many bits to interpret)
- Overflow page pointer size – 1 byte (number of bytes to read to interpret the overflow page number)
- Record entry size – 1 byte (the size of each record entry)
- Bytes reserved to store number of record entries are stored in index page (exclude directory page) – 1 byte (number of bytes to read to interpret the number of records in the current page)
- Directory entry size (size of each directory entry)
- A directory of buckets that points to the index page for the rest of the header page (a list of fixed-size directory entries)
- Overflow pointer for the next directory page – 2 bytes (the page number that points to the next directory page)



Linear Hash Index File **directory page** Format contains the following information in order:

- Number of directory entries in this page – 2 bytes (number of fixed size directory entries in this page)

- A list of directory entries (a sequential list of fixed size directory entries)
- Overflow directory pointer to the next directory page – 2 bytes (the page number for the next directory page)



Linear Hash Index File index and overflow page Format:

- Number of records in this page – 2 bytes (number of fixed size record entries in this page)
- All the record entries (a sequential list of fixed size record entries)
- At the end, there is an overflow page pointer – 2 bytes (the page number for the overflow page of this bucket)



```

75 -has overflow pointer
76 ===
77 def __init__(self, size, pagellum, entrySize=15):
78     ExtendablePage.__init__(self, size, pagellum)
79
80 if __name__ == "__main__":
81     pSize = 1024
82     field = 0
83     indexType = 1
84     heapFilePath = "../names.db"
85     numBuckets = 512
86
87     import os
88     byteSize = os.path.getsize(heapFilePath)
89     # Construct a in-memory index file with proper format
90     staticIndex = staticHashIndexer("staticIndex", numBuckets, pSize, field, indexType, DATAENTRYSIZE)
91     with open(heapFilePath, "rb") as heapfile:
92         rowid = 0
93         while True:
94             heapfile.seek(RECORDSIZE * rowid, 0)
95             if heapfile.tell() >= byteSize:
96                 break
97             # Read the records with null characters behind
98             firstName = heapfile.read(FIRSTNAMELENGTH).decode("utf-8")
99             staticIndex.hash(firstName, rowid)
100             rowid += 1
101
102     # write the index file to disk
103     staticIndex.writeToDisk()

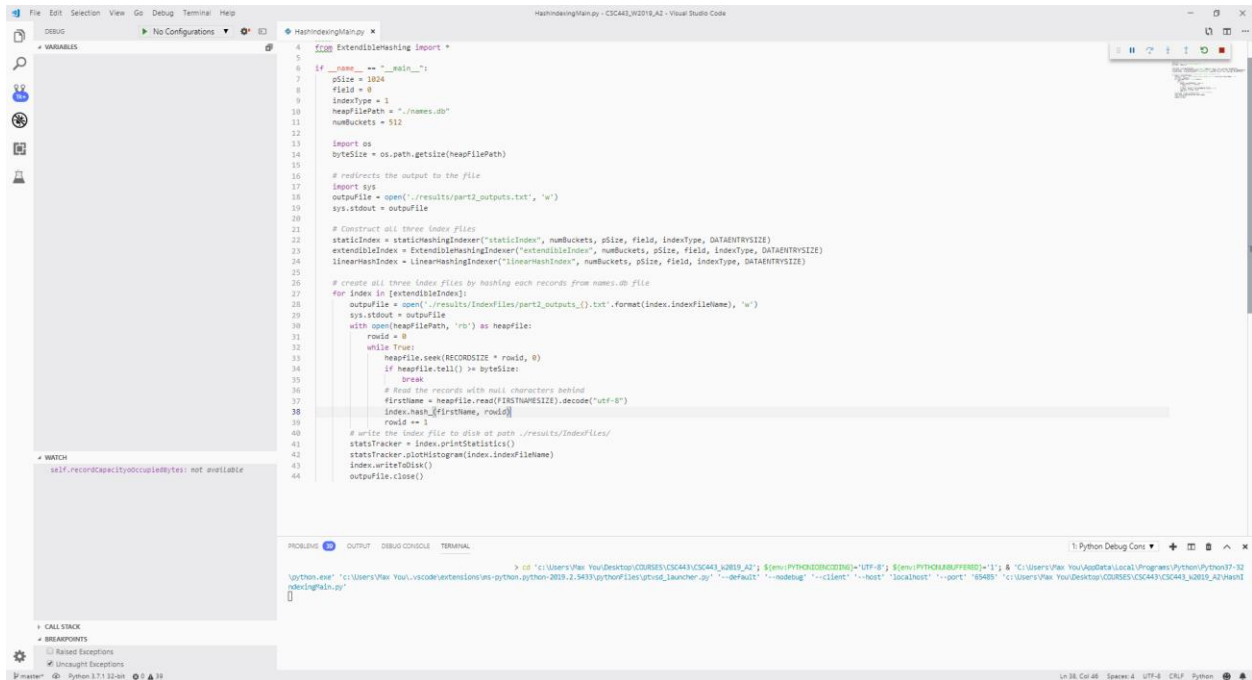
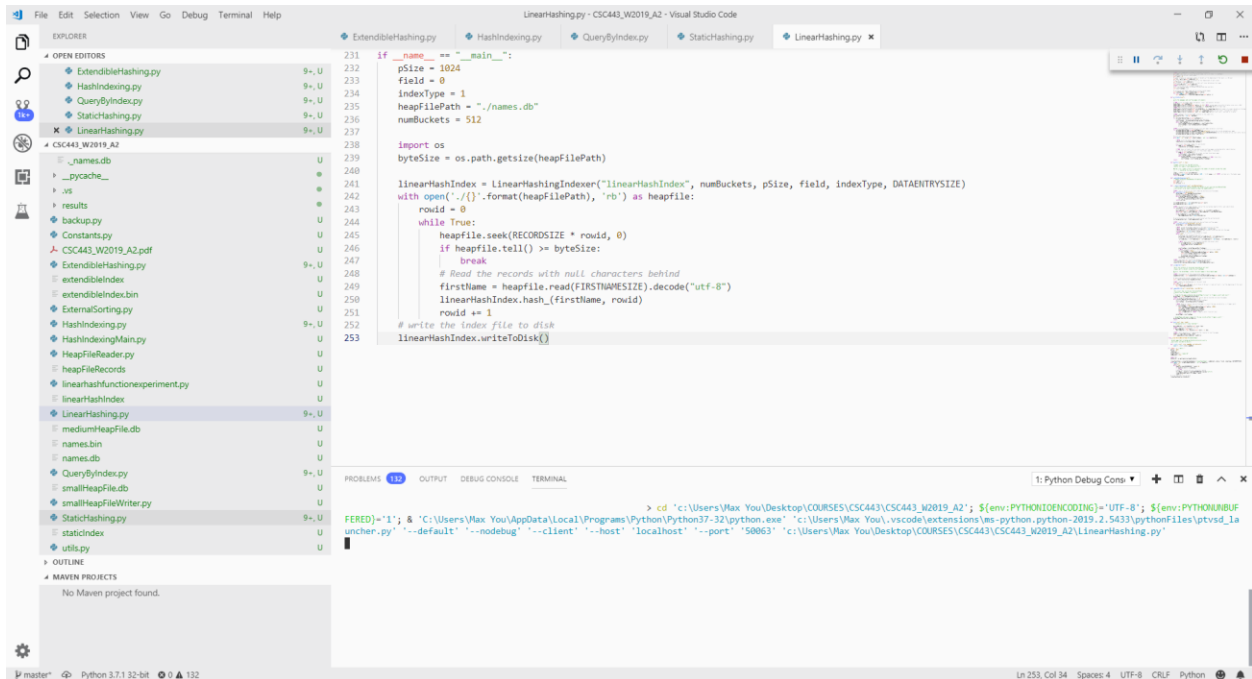
```

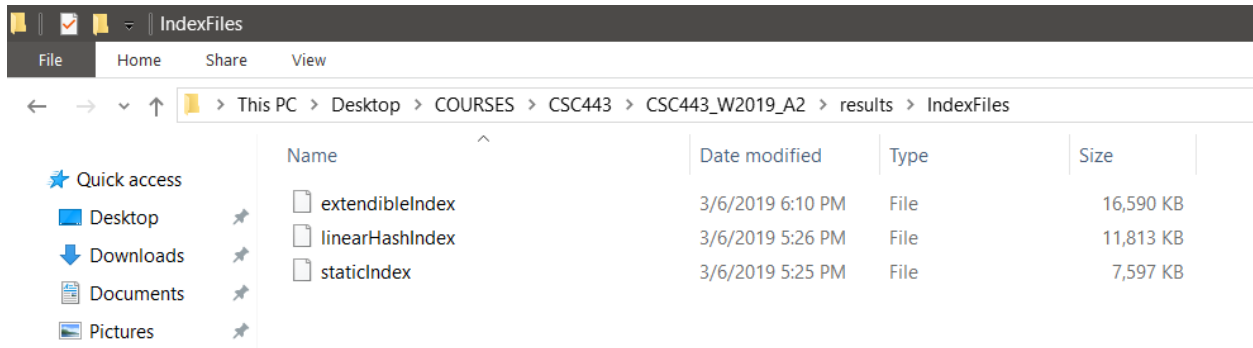
Terminal output:

```

> cd 'c:\Users\Max You\Desktop\COURSES\CSC443\CSC443_W2019_A2'; $([env:PYTHONIDENCODING]=UTF-8; $([env:PYTHONUNBUFFERED]=1; & 'c:\Users\Max You\AppData\Local\Programs\Python\Python37-32\python.exe' 'c:\Users\Max You\.vscode\extensions\ms-python.python-2019.2.5433\pythonFiles\ptvsd_launcher.py' '--default' '--nodebug' '--client' '--host' 'localhost' '--port' '50857' 'c:\Users\Max You\Desktop\COURSES\CSC443\CSC443_W2019_A2\StaticHashing.py')

```



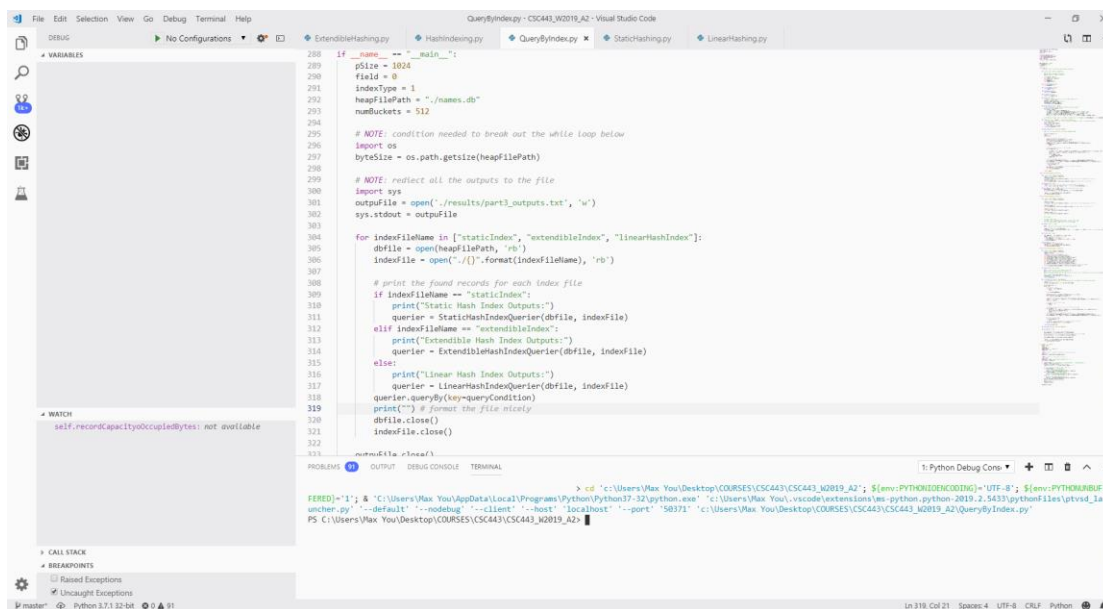


Part three answers:

Index type \ metadata	Index Page Reads	Data Page Reads	Index Page Writes	Data Page Writes
Static hashing index	18	39	0	0
Extendible hashing index	3	39	0	0
Linear hashing index	2	39	0	0

Observations:

Firstly, we noticed that static hashing index has the most index page reads since mostly like it has the greatest number of overflow pages at that bucket due to its property of fixed number of buckets. Secondly, we see that extendible hash index has 3 index page reads, compare to other hash index files, it needs to read an extra page that contains the desired directory entry. Lastly, we noticed that linear hash index has very little index page reads, it makes sense when combine the histogram results from above, we see most of the buckets have only 1-2 index pages. The number of data page reads should be the same across all three type of index file since record entry contains the rowid as the value; each time we see the rowid with the matched first name, we need to fetch a page from the database file to get the actual records. Finally, query operations should not have write cost, we only need to read in the data page to find the records.



The queried results are the same across all three hash index files, as follows:

First name: Nona - Last name: Buentello - Email: nona.buentello@yahoo.com

First name: Nona - Last name: Kluesner - Email: nona.kluesner@hotmail.com

First name: Nona - Last name: Burgoon - Email: nona.burgoon@aol.com

First name: Nona - Last name: Bloomberg - Email: nona.bloomberg@gmail.com

First name: Nona - Last name: Gilmartin - Email: nona.gilmartin@gmail.com

First name: Nona - Last name: Evan - Email: nona.evan@yahoo.co.uk

First name: Nona - Last name: Rowlands - Email: nona.rowlands@aol.com

First name: Nona - Last name: Cabe - Email: nona.cabe@aol.com

First name: Nona - Last name: Quackenbush - Email: nona.quackenbush@charter.net

First name: Nona - Last name: Lucht - Email: nona.lucht@yahoo.ca

First name: Nona - Last name: Mapp - Email: nona.mapp@aol.com

First name: Nona - Last name: Enlow - Email: nona.enlow@microsoft.com

First name: Nona - Last name: Hulbert - Email: nona.hulbert@msn.com

First name: Nona - Last name: Mcclurg - Email: nona.mcclurg@rediffmail.com

First name: Nona - Last name: Neil - Email: nona.neil@msn.com

First name: Nona - Last name: Sipp - Email: nona.sipp@gmail.com

First name: Nona - Last name: Catchings - Email: nona.catchings@bellsouth.net

First name: Nona - Last name: Skalski - Email: nona.skalski@gmail.com

First name: Nona - Last name: McCloskey - Email: nona.mccloskey@aol.com

First name: Nona - Last name: Lilly - Email: nona.lilly@yahoo.com

First name: Nona - Last name: Hirth - Email: nona.hirth@verizon.net

First name: Nona - Last name: Arias - Email: nona.arias@aol.com

First name: Nona - Last name: Lansberry - Email: nona.lansberry@yahoo.com

First name: Nona - Last name: Cote - Email: nona.cote@yahoo.co.in

First name: Nona - Last name: Forest - Email: nona.forest@exxonmobil.com

First name: Nona - Last name: Vizcaino - Email: nona.vizcaino@gmail.com

First name: Nona - Last name: Achenbach - Email: nona.achenbach@rediffmail.com

First name: Nona - Last name: Belliveau - Email: nona.belliveau@gmail.com

First name: Nona - Last name: Sheehan - Email: nona.sheehan@hotmail.com

First name: Nona - Last name: Edgar - Email: nona.edgar@yahoo.com

First name: Nona - Last name: Cyr - Email: nona.cyr@gmail.com

First name: Nona - Last name: Peek - Email: nona.peek@shell.com

First name: Nona - Last name: Budd - Email: nona.budd@gmail.com

First name: Nona - Last name: Preston - Email: nona.preston@aol.com

First name: Nona - Last name: Vizcaino - Email: nona.vizcaino@gmail.com

First name: Nona - Last name: Snook - Email: nona.snook@hotmail.com

First name: Nona - Last name: Lamas - Email: nona.lamas@apple.com

First name: Nona - Last name: Ventimiglia - Email: nona.ventimiglia@gmail.com

First name: Nona - Last name: Woodworth - Email: nona.woodworth@gmail.com