

Geographic Information System

Geographic information systems organize information pertaining to geographic features and provide various kinds of access to the information. A geographic feature may possess many attributes (see below). In particular, a geographic feature has a specific location. There are a number of ways to specify location. For this project, we will use latitude and longitude, which will allow us to deal with geographic features at any location on earth. A reasonably detailed tutorial on latitude and longitude can be found in the Wikipedia at en.wikipedia.org/wiki/Latitude and en.wikipedia.org/wiki/Longitude.

The GIS record files were obtained from the website for the USGS Board on Geographic Names (geonames.usgs.gov). The file begins with a descriptive header line, followed by a sequence of GIS records, one per line, which contain the following fields in the indicated order (all are mandatory unless indicated otherwise):

Figure 1: Geographic Data Record Format

Significance	Type/Format	Comments
Feature ID (FID)	non-negative integer	unique identifier for this geographic feature
Feature name	string	standard name of feature
Feature class	string	descriptive classification of feature
State alphabetic code	two-characters	US postal code abbreviation
State numeric code	non-negative integer	numeric code for state
County name	string	county in which feature occurs
County numeric code	non-negative integer	numeric code for county
Primary latitude (DMS)	DDMMSS['N' 'S']	feature latitude in DMS format or Unknown
Primary longitude (DMS)	DDMMSS['E' 'W']	feature longitude in DMS format or Unknown
Primary latitude (dec deg)	decimal number	feature latitude in decimal format or Unknown
Primary longitude (dec deg)	decimal number	feature longitude in decimal format or Unknown
Source latitude (DMS)	DDMMSS['N' 'S']	latitude of feature source in DMS format, optional
Source longitude (DMS)	DDMMSS['E' 'W']	longitude of feature source in DMS format, optional
Source latitude (dec deg)	decimal number	latitude of feature source in decimal format, optional
Source longitude (dec deg)	decimal number	longitude of feature source in decimal format, optional
Feature elevation in meters	integer	altitude above/below sea level, optional
Feature elevation in feet	integer	altitude above/below sea level, optional
Map name	string	name of USGS topographic map including feature
Date created	string	date feature was initially committed to the database
Date edited	string	date feature record was last updated, optional

In the GIS record file, each record will occur on a single line, and the fields will be separated by pipe (`|`) symbols. Empty fields will be indicated by a pair of pipe symbols with no characters between them. See the posted `VA_Monterey.txt` file for many examples.

GIS record files are guaranteed to conform to this syntax, so there is no explicit requirement that you validate the files. On the other hand, some error-checking during parsing may help you detect errors in your parsing logic.

The file can be thought of as a sequence of bytes, each at a unique offset from the beginning of the file, just like the cells of an array. So, each GIS record begins at a unique offset from the beginning of the file.

Note:

Each line of a text file ends with a particular marker (known as the line terminator). In MS-DOS/Windows file systems, the line terminator is a sequence of two ASCII characters (CR + LF). In Unix systems, the line terminator is a single ASCII character (LF). Other systems may use other line termination conventions.

Why should you care? Which line termination is used has an effect on the file offsets for all but the first record in the data file. As long as we're all testing with files that use the same line termination, we should all get the same file offsets. But if you change the file format (of the posted data files) to use different line termination, you will get different file offsets than are shown in the posted log files. Most good text editors will tell you what line termination is used in an opened file, and also let you change the line termination scheme.

Figure 2: Sample Geographic Data Records

Note that some record fields are optional, and that when there is no given value for a field, there are still delimiter symbols for it.

Also, some of the lines are "wrapped" to fit into the text box; lines are never "wrapped" in the actual data files.

```

FEATURE_ID|FEATURE_NAME|FEATURE_CLASS|STATE_ALPHA|STATE_NUMERIC|COUNTY_NAME|COUNTY_NUMERIC|PRIMARY_LAT_DMS|PRIM_LONG_DMS|PRIM_LAT_DEC|PRIM_LONG_DEC|SOURCE_LAT_DMS|SOURCE_LONG_DMS|SOURCE_LAT_DEC|SOURCE_LONG_DEC|ELEV_IN_FT|MAP_NAME|DATE_CREATED|DATE_EDITED
1479116|Monterey Elementary School|School|VA|51|Roanoke (city)|770|371906N|0795608W|37.3183753|-
79.9355857|111|323|1060|Roanoke|09/28/1979|09/15/2010
1481345|Asbury Church|Church|VA|51|Highland|091|382607N|0793312W|38.4353981|-79.5533807|111|1818|2684|Monterey|09/28/1979|
1481852|Blue Grass Populated Place|VA|51|Highland|091|383000N|0793259W|38.5001188|-79.5497702|111|1777|2549|Monterey|09/28/1979|
1481878|Bluegrass Valley|Valley|VA|51|Highland|091|382953N|0793322W|38.4981745|-79.5394921|382601N|0793800W|38.4337309|-
79.6333833|759|2490|Monterey|09/28/1979|
1482110|Buck Hill|Summit|VA|51|Highland|091|381902N|0793358W|38.3173452|-79.5661577|111|1003|3291|Monterey SE|09/28/1979|
1482176|Burners Run|Stream|VA|51|Highland|091|382509N|0793409W|38.4192873|-79.5692144|382531N|0793538W|38.4252778|-
79.5938889|848|2782|Monterey|09/28/1979|
1482324|Mount Carlyle|Summit|VA|51|Highland|091|381556N|0793353W|38.2656799|-79.5647682|111|1698|2290|Monterey SE|09/28/1979|
1482434|Central Church|Church|VA|51|Highland|091|382953N|0793323W|38.4981744|-79.5564371|111|1773|2536|Monterey|09/28/1979|
1482557|Claylick Hollow|Valley|VA|51|Highland|091|381613N|0793238W|38.2704021|-79.5439343|381733N|0793324W|38.29251|-
79.5566667|573|1880|Monterey SE|09/28/1979|
1482785|Crab Run|Stream|VA|51|Highland|091|381707N|0793144W|38.2854018|-79.528934|381903N|0793415W|38.3175|-79.5708333|579|1900|Monterey SE|09/28/1979|
1482950|Davis Run|Stream|VA|51|Highland|091|381824N|0793053W|38.3067903|-79.5147671|382057N|0793505W|38.3491667|-79.5847222|601|1972|Monterey SE|09/28/1979|
1483281|Elk Run|Stream|VA|51|Highland|091|382936N|0793153W|38.4934524|-79.5314362|383121N|0793056W|38.5226185|-
79.5156027|757|2484|Monterey|09/28/1979|
1483492|Forks of Waters|Locale|VA|51|Highland|091|382856N|0793031W|38.4823417|-79.5086575|111|1705|2313|Monterey|09/28/1979|
1483527|Frank Run|Stream|VA|51|Highland|091|382953N|0793310W|38.4981744|-79.5528258|383304N|0793341W|38.5512285|-
79.5614381|780|2559|Monterey|09/28/1979|
1483647|Ginseng Mountain|Summit|VA|51|Highland|091|382850N|0793139W|38.4806751|-79.527547|111|1978|3209|Monterey|09/28/1979|
1483860|Gulf Mountain|Summit|VA|51|Highland|091|382940N|0793103W|38.4945636|-79.5175468|111|1006|3300|Monterey|09/28/1979|
1483916|Hamilton Chapel|Church|VA|51|Highland|091|381740N|0793707W|38.2945677|-79.6186591|111|1823|2700|Monterey SE|09/28/1979|
1484097|Highland High School|School|VA|51|Highland|091|382426N|0793444W|38.4071387|-79.5789333|111|1879|2884|Monterey|09/28/1979|09/15/2010
1484099|Highland Wildlife Management Area|Park|VA|51|Highland|091|381905N|0793439W|38.3181785|-79.5775471|111|1954|3130|Monterey SE|09/28/1979|
. . .
    
```

Assignment:

You will implement a system that indexes and provides search features for a file of GIS records, as described above.

Your system will build and maintain several in-memory index data structures to support these operations:

- Importing new GIS records into the database file
- Retrieving data for all GIS records matching given geographic coordinates
- Retrieving data for all GIS records matching a given feature name and state
- Retrieving data for all GIS records that fall within a given (rectangular) geographic region
- Displaying the in-memory indices in a human-readable manner

You will implement a single software system in Java to perform all system functions.

Program Invocation:

The program will take the names of three files from the command line, like this:

```
java GIS <database file name> <command script file name> <log file name>
```

The database file should be created as an empty file; note that the specified database file may already exist, in which case the existing file should be truncated or deleted and recreated. If the command script file is not found the program should write an error message to the console and exit. The log file should be rewritten every time the program is run, so if the file already exists it should be truncated or deleted and recreated.

System Overview:

The system will create and maintain a GIS database file that contains all the records that are imported as the program runs. The GIS database file will be empty initially. All the indexing of records will be done relative to this file.

There is no guarantee that the GIS record file will not contain two or more distinct records that have the same geographic coordinates. In fact, this is natural since the coordinates are expressed in the usual DMS system. So, we cannot treat geographic coordinates as a primary (unique) key.

The GIS records will be indexed by the **Feature Name** and **State** (abbreviation) fields. This *name index* will support finding offsets of GIS records that match a given feature name and state abbreviation.

The GIS records will also be indexed by geographic coordinate. This *coordinate index* will support finding offsets of GIS records that match a given primary latitude and primary longitude.

The system will include a *buffer pool*, as a front end for the GIS database file, to improve search speed. See the discussion of the buffer pool below for detailed requirements. When performing searches, retrieving a GIS record from the database file must be managed through the buffer pool. During an import operation, when records are written to the database file, the buffer pool will be bypassed, since the buffer pool would not improve performance during imports.

When searches are performed, complete GIS records will be retrieved from the GIS database file that your program maintains. The only complete GIS records that are stored in memory at any time are those that have just been retrieved to satisfy the current search, or individual GIS records created while importing data or GIS records stored in the buffer pool.

Aside from where specific data structures are required, you may use any suitable Java library containers you like.

Each index should have the ability to supply a nicely-formatted representation of itself as a `String` object, or to write a nicely-formatted display of itself to an output stream.

Name Index Internals:

The *name index* will use a hash table for its physical organization. Each hash table entry will store a feature name and state abbreviation (separately or concatenated, as you like) and the file offset(s) of the matching record(s). Since each GIS record occupies one line in the file, it is a trivial matter to locate and read a record given nothing but the file offset at which the record begins.

The hash table must use a contiguous physical structure (array). The initial size of the table will be 1019, and the table will resize itself automatically. The precise logic for resizing is up to you; if you want to match my output logs, increase the size of the table to the next value in the list below when the table becomes 70% full:

1019, 2027, 4079, 8123, 16267, 32503, 65011, 130027, 260111, 520279, 1040387, 2080763, 4161539, 8323151, 16646323

Your table will use quadratic probing to resolve collisions, with the quadratic function $(n^2 + n)/2$ to compute the step size. Since the table sizes given above are all primes of the form $4k + 3$, an empty slot will always be found unless the table is full.

You will use the `elfhash()` function from the course notes, and apply it to the concatenation of the feature name and state (abbreviation) field of the data records. Precisely how you form the concatenation is up to you.

You must be able to display the contents of the hash table in a readable manner.

Coordinate Index Internals:

The coordinate index will use a *bucket* PR quadtree for the physical organization. In a bucket PR quadtree, each leaf stores up to K data objects (for some fixed value of K). Upon insertion, if the added value would fall into a leaf that is already full, then the region corresponding to the leaf will be partitioned into quadrants and the $K+1$ data objects will be inserted into those quadrants as appropriate. As is the case with the regular PR quadtree, this may lead to a sequence of partitioning steps, extending the relevant branch of the quadtree by multiple levels. In this project, K will probably equal 4, but I reserve the right to specify a different bucket size with little notice, so this should be easy to modify.

The index entries held in the quadtree will store a geographic coordinate and a collection of the file offsets of the matching GIS records in the database file.

Note: do not confuse the bucket size with any limit on the number of GIS records that may be associated with a single geographic coordinate. A quadtree node can contain index objects for up to K different geographic coordinates. Each such index object can contain references to an unlimited number of different GIS records.

The PR quadtree implementation should follow good design practices, and its interface should be somewhat similar to that of the BST. You are expected to implement different types for the leaf and internal nodes, with appropriate data membership for each, and an abstract base type from which they are both derived. Of course, these were all requirements for the related minor project.

You must be able to display the PR quadtree in a readable manner. PR quadtree display code is given in the course notes. The display must clearly indicate the structure of the tree, the relationships between its nodes, and the data objects in the leaf nodes.

Buffer Pool Details:

The buffer pool for the database file should be capable of buffering up to 20 records, and will use LRU replacement. You may use any structure you like to organize the pool slots; however, since the pool will have to deal with record replacements, some structures will be more efficient (and simpler) to use.

It is up to you to decide whether your buffer pool stores interpreted or raw data; i.e., whether the buffer pool stores GIS record objects or just strings.

You must be able to display the contents of the buffer pool, listed from MRU to LRU entry, in a readable manner. The order in which you retrieve records when servicing a multi-match search is not specified, so such searches may result in different orderings of the records within the buffer pool. That is OK.

A Note on Coordinates and Spatial Regions:

It is important to remember that there are fundamental differences between the notion that a geographic feature has specific coordinates (which may be thought of as a point) and the notion that each node of the PR quadtree corresponds to a particular sub-region of the coordinate space (which will usually contain many points).

In this assignment, coordinates of geographic features are specified as latitude/longitude pairs, and the minimum resolution is one second of arc. Thus, you may think of the geographic coordinates as being specified by a pair of integer values.

On the other hand, the boundaries of the sub-regions are determined by performing arithmetic operations, including division, starting with the values that define the boundaries of the “world”. Unless the dimensions of the world happen to be powers of 2, this can quickly lead to regions whose boundaries cannot be expressed exactly as integer values. You may use floating-point values or integer values to represent region boundaries when computing region boundaries during splitting and quadtree traversals. If you use integers, be careful not to unintentionally create “gaps” between regions.

Your implementation should view the boundary between regions as belonging to one of those regions. The choice of a particular rule for handling this situation is left to you. The specification for the minor PR quadtree project describes how I made that decision, but there is absolutely no requirement that you follow the same approach.

When carrying out a region search, you must determine whether the search region overlaps with the region corresponding to a subtree node before descending into that subtree. The Java libraries include a `Rectangle` class which could be (too) useful. You may make use of the `Rectangle` class, but you will be penalized 10% if your submitted solution makes use of any of the following `Rectangle` methods: `contains()`, `intersection()`, and `intersects()`. Note though, that it is acceptable to make use of those methods during development, but you must implement your own versions of them in your final submission.

Finally, do not interpret longitude/latitude values sloppily. In particular, do not represent a value given in DMS format, like 1214322W, as the integer -1214322. (The negative sign comes from the fact that it's west longitude, and that's fine.) Some students have encountered problems in the past when using that approach. Instead, convert the given value to total seconds; in this case that would be -438202 seconds.

Other System Elements:

There should be an overall controller that validates the command line arguments and manages the initialization of the various system components. The controller should hand off execution to a command processor that manages retrieving commands from the script file, and making the necessary calls to other components in order to carry out those commands.

Naturally, there should be a data type that models a GIS record.

There may well be additional system elements, whether data types or data structures, or system components that are not mentioned here. The fact no additional elements are explicitly identified here does not imply that you will not be expected to analyze the design issues carefully, and to perhaps include such elements.

Aside from the command-line interface, there are no specific requirements for interfaces of any of the classes that will make up your GIS; it is up to you to analyze the specification and come up with an appropriate set of classes, and to design their interfaces to facilitate the necessary interactions. It is probably worth pointing out that an index (e.g., a geographic coordinate index) should not simply be a naked container object (e.g. quadtree); if that's not clear to you, think more carefully about what sort of interface would be appropriate for an index, as opposed to a container.

Command File:

The execution of the program will be driven by a script file. Lines beginning with a semicolon character (';') are comments and should be ignored. Each non-comment line of the command file will specify one of the commands described below.

Each line consists of a sequence of tokens, which will be separated by single tab characters. A newline character will immediately follow the final token on each line. The command file is guaranteed to conform to this specification, so you do not need to worry about error-checking when reading it. The following commands must be supported:

`world<tab><westLong><tab><eastLong><tab><southLat><tab><northLat>`

This will be the first command in the file, and will occur once. It specifies the boundaries of the coordinate space to be modeled. The four parameters will be longitude and latitudes expressed in DMS format, representing the vertical and horizontal boundaries of the coordinate space.

It is possible that the GIS record file will contain records for features that lie outside the specified coordinate space. Such records should be ignored; i.e., they will not be indexed.

`import<tab><GIS record file>`

Add all the GIS records in the specified file to the database file. This means that the records will be appended to the existing database file, and that those records will be indexed in the manner described earlier. When the import is completed, log the number of entries added to each index, and the longest probe sequence that was needed when inserting to the hash table.

`what_is_at<tab><geographic coordinate>`

For every GIS record in the database file that matches the given `<geographic coordinate>`, log the offset at which the record was found, and the feature name, county name, and state abbreviation. Do not log any other data from the records.

`what_is<tab><feature name><tab><state abbreviation>`

For every GIS record in the database file that matches the given `<feature name>` and `<state abbreviation>`, log the offset at which the record was found, and the county name, the primary latitude, and the primary longitude. Do not log any other data from the records.

`what_is_in<tab><geographic coordinate><tab><half-height><tab><half-width>`

For every GIS record in the database file whose coordinates fall within the closed rectangle with the specified height and width, centered at the `<geographic coordinate>`, log the offset at which the record was found, and the feature name, the state name, and the primary latitude and primary longitude. Do not log any other data from the records. The half-height and half-width are specified as seconds.

`what_is_in<tab>-l<tab><geographic coordinate><tab><half-height><tab><half-width>`

For every GIS record in the database file whose coordinates fall within the closed rectangle with the specified height and width, centered at the `<geographic coordinate>`, log every important non-empty field, nicely formatted and labeled. See the posted log files for an example. Do not log any empty fields. The half-height and half-width are specified as seconds.

`what_is_in<tab>-c<tab><geographic coordinate><tab><half-height><tab><half-width>`

Log the number of GIS records in the database file whose coordinates fall within the closed rectangle with the specified height and width, centered at the `<geographic coordinate>`. Do not log any data from the records themselves. The half-height and half-width are specified as seconds.

`debug<tab>[quad | hash | pool]`

Log the contents of the specified index structure in a fashion that makes the internal structure and contents of the index clear. It is not necessary to be overly verbose here, but it would be useful to include information like key values and file offsets where appropriate.

`quit<tab>`

Terminate program execution.

If a `<geographic coordinate>` is specified for a command, it will be expressed as a pair of latitude/longitude values, expressed in the same DMS format that is used in the GIS record files.

For all the commands, if a search results in displaying information about multiple records, you may display that data in any order that is natural to your design.

Sample command scripts will be provided on the website. As a general rule, every command should result in some output. In particular, a descriptive message should be logged if a search yields no matching records.

Log File Description:

Since this assignment will be graded by TAs, rather than the Curator, the format of the output is left up to you. Of course, your output should be clear, concise, well labeled, and correct. You should begin the log with a few lines identifying yourself, and listing the names of the input files that are being used.

The remainder of the log file output should come directly from your processing of the command file. You are required to echo each comment line, and each command that you process to the log file so that it's easy to determine which command each section of your output corresponds to. Each command (except for "world") should be numbered, starting with 1, and the output from each command should be well formatted, and delimited from the output resulting from processing other commands. A complete sample log will be posted shortly on the course website.

Administrative Issues:

Submission:

You will submit this assignment to the Curator System (read the *Student Guide*), where it will be archived for grading at a demo with a TA.

For this assignment, you must submit a `jar` file containing all the Java source code files for your implementation (i.e., `.java` files). Submit only the Java source files. Do not submit Java bytecode (class) files. Submit nothing else.

In order to correct submission errors and late-breaking implementation errors, you will be allowed up to five submissions for this assignment. You may choose which one will be evaluated at your demo.

The Student Guide and link to the submission client can be found at: <http://www.cs.vt.edu/curator/>

Evaluation:

The quality of the OO design in your solution will be evaluated; you are expected to identify and implement a sound collection of classes, whether the specification mentions them or not. The quality of your internal documentation will be evaluated. Finally, the correctness of your solution will be evaluated by executing your solution on a collection of test data files. Be sure to test your solution with all of the data sets that will be posted, since we will use a variety of data sets, including at least one very large data one (perhaps hundreds of thousands of records) in our evaluation.

The quality of your design will determine a substantial portion of your grade for this assignment. Pay attention to the discussion and solution for the related design assignment, Homework 0, and make changes as needed. It is possible you may lose points during the project evaluation for poor design choices, even if you have already lost points for the same poor decisions on Homework 0. On the other hand, if you lose points on Homework 0, but your implementation of this project indicates that you have corrected those decisions in your final design, we will rebate part of the deductions you received on Homework 0.

Remember that your implementation will be tested in the specified environment. If you use a different development platform, it is entirely your responsibility to make sure your implementation works correctly in the lab.

Note that the evaluation of your project will depend substantially on the quality of your code and documentation. See the Programming Standards page on the course website for specific requirements that should be observed in this course.

Pedagogic points:

The goals of this assignment include, but are not limited to:

- implementation of a hash table generic in Java
- implementation of a bucket PR quadtree generic in Java
- implementation of a buffer pool in Java
- implementation of complementary internal and leaf node types to conserve memory, using an appropriate hierarchy of types
- design of appropriate index classes to wrap the containers
- design of appropriate data objects to store in each index
- understanding how to navigate a file in Java
- creation of a sensible OO design for the overall system, including the identification of a number of useful classes not explicitly named in this specification
- implementation of such an OO design into a working system
- incremental testing of the basic components of the system in isolation
- satisfaction when the entire system comes together in good working order

Pledge:

Each of your program submissions must be pledged to conform to the Honor Code requirements for this course. Specifically, you **must** include the pledge statement provided on the Projects page of the course website. The pledge statement should be in the file that contains your main class.