# Project 2 Document

Author: Yan Wang & Zishan Qin

## Project Requirement

1. An indexing mechanism:B+ tree or Hash function index..
2. Five interfaces related to: Put, Get and Remove.
3. Extra issues: multiple threadings, value-based search

## Assumption

The total size of value store will not exceed the limit of memory

## Design

1. We built a class called Btree for building up our index mechanism and doing Put/Get/Remove operations. Btree is composed by nodes, each node has a list of entrys. Btree class is used for maintaining the tree size, tree height and root information. Node is used for maintaing the usage of entries in this node. Entry is special : for internal nodes,only key and next pointer will be used; for external nodes, only key and value will be used.

2. The dataset used in this project is a list of words for indexing the bag-of-word data(Supported by UCI datasets). It's a Key-Value pair dataset and the size of tuples is more than 32000. So we can make the fanout  of b+ tree as 10. So for get/put/remove operations, only 6 times look-ups will be needed.

3. Although the description of the project only ask for two kinds of Get/Put/Remove Key-Value pair data types(`void Put(string key, Number data_value);` or`void Put(string key, string data_value);`). We implemented using a more general way- using Java Generics to abstract the Key-Value pair,as BTree<Key extends Comparable<Key>, Value>.   Key and Value can be any object type(Integer/String/Date...) as long as the Key is a comparable object type.  So we only need to implement 3 functions but the codes are applicable for the five interfaces in the project description.

   For Put/Get/Remove functions, we implemented in this way, supposed the fanout is m:

   Get: we start from root, based on comparison of keys of upper layer, we step into next layer untill we get height=0.When height=0, we are in the Leaf-layer, we locate the item according to its key. The time complexity is $\log_m(n)$.

   Put:we start from root, locate the biggest leafnode which is smaller than the inserted item, we try to insert it in the entries of that leaf node. If the entries size is less than M, then we are done. If node, we need to split that node into two and insert the smaller

spliited node into its parent entries.If the parent is full, split it too, repeat the split process above until a parent is found that need not split.

Remove:Perform a search to determine which leaf node contains the key. Remove the key from the leaf node. If the leaf node is at least half-full, done! If the leaf node as L is less than half-full: Try to borrow a key from sibling node as S (adjacent node with same parent). If S is L's left sibling, then borrow S's last key, and replace their parent navigate key with this borrowed key value. If S is L's right sibling, then borrow S's first key, and replace their parent navigate key with S's second key value. If can not borrow a key from sibling node, then merge L and sibling S
After merged L and S, delete their parent navigate key and proper child pointer. Repeat the borrow or merge operation on parent node, perhaps propagate to root node and decrease the height of the tree.

4. All three methods "put", "get", "remove" are synchronized. Every method will wait until it is notified by other threadings. After the method is implemented, it will notify other threadings.
5. For value-based search, since we use Key to build B+ tree so in this context, so B+ tree is not helpful,  so we implemented a new function  Key get(Value val) to do the search.

# Test

**Test indexing mechanism:**
we read data from the dataset line by line and simutaneously, insert these data into btree. After that we check the properties of that Btree, the height is 6 and there are 32108 records stored.

```
############ Test B+ tree Build-up #######################
B+ tree with height of 6
B+ tree with records of 32108
```

**Test Five interfaces related to: Put, Get and Remove.**
in this test we do folloing things:
   1) Get(key=96886)
   2) Insert(1000,a3333.txt) into btree.
   3) Get(key=1000);
   4) Check the tree size: notice there the size of tree is increased by one.
   5) Remove(Key=1000);
   6) check the tree size:notice there the size of tree is decreased by one.
   7) Get(key=1000)

```
########### tree B+ tree operations #######################
Test B+tree Get
Get record with key 96886: Record (96886)-a0000536
Test B+tree Put
Get record with key 1000: Record (1000)-a33333.txt
B+ tree with records of 32109
Test B+tree Remove with key 1000
B+ tree with records of 32108
Try to get record with key 1000...  Not found!
```

**Test multiple threadings**

Multiple threadings are generated to mimic multiple client. Except for the first threading, we don't control the order of other threadings. Therefore, we generated a bunch of threadings for "get", so as to display the "put" and "remove" go well. Result of this test shows below.

```
############ Test Multithreads #############################
ThreadGet  get value a0000009 with key 96697
ThreadGet  get value a0000009 with key 96697
ThreadGet  get value a0000009 with key 96697
ThreadGet  get value a0000009 with key 96697
ThreadRemove  remove the item with key 96697
ThreadGet try to get value with key 96697, but not found
                                                      ...
```

**Test Value-Based Search**

For value-based search, since we use Key to build B+ tree so in this context, so B+ tree is not helpful,  so we implemented a new function  Key get(Value val) to do the search. we randomly select a record(100225,"a0074797"), and do the scanning on Value-"a0074797" line by line and finally find the key 100225.

```
############ Test Get Function based on Value #############
Test Get(Value value) function
Get record with Value (a0074797): Key-100225
```