

Join GitHub today

Dismiss

GitHub is home to over 31 million developers working together to host and review code, manage projects, and build software together.

[Sign up](#)

Branch: master ▾ Fa18HW5 / src / main / java / edu / berkeley / cs186 / database / index / InnerNode.java

[Find file](#) [Copy path](#)

es1024 HW5 skeleton code

16ecdac on Nov 8, 2018

1 contributor

463 lines (415 sloc) 16.4 KB

[Raw](#) [Blame](#) [History](#)   

```
1 package edu.berkeley.cs186.database.index;
2
3 import java.nio.ByteBuffer;
4 import java.util.*;
5
6 import edu.berkeley.cs186.database.BaseTransaction;
7 import edu.berkeley.cs186.database.common.Buffer;
8 import edu.berkeley.cs186.database.common.Pair;
9 import edu.berkeley.cs186.database.databox.DataBox;
10 import edu.berkeley.cs186.database.databox.Type;
11 import edu.berkeley.cs186.database.io.Page;
12 import edu.berkeley.cs186.database.table.RecordId;
13
14 /**
15  * A inner node of a B+ tree. Every inner node in a B+ tree of order d stores
16  * between d and 2d keys. An inner node with n keys stores n + 1 "pointers" to
17  * children nodes (where a pointer is just a page number). Moreover, every
18  * inner node is serialized and persisted on a single page; see toBytes and
19  * fromBytes for details on how an inner node is serialized. For example, here
20  * is an illustration of an order 2 inner node:
21  *
22  * +---+---+---+---+
23  * | 10 | 20 | 30 |   |
24  * +---+---+---+---+
25  * /   |   |   \
26  */
27 class InnerNode extends BPlusNode {
28     // Metadata about the B+ tree that this node belongs to.
29     private BPlusTreeMetadata metadata;
30
31     // The page on which this leaf is serialized.
32     private Page page;
33
34     // The keys and child pointers of this inner node. See the comment above
35     // LeafNode.keys and LeafNode.rids in LeafNode.java for a warning on the
36     // difference between the keys and children here versus the keys and children
37     // stored on disk.
38     private List<DataBox> keys;
39     private List<Integer> children;
40
41     // Constructors //////////////////////////////////////
42     /**
```

```

43 * Construct a brand new inner node. The inner node will be persisted on a
44 * brand new page allocated by metadata.getAllocator().
45 */
46 public InnerNode(BPlusTreeMetadata metadata, List<DataBox> keys,
47                 List<Integer> children, BaseTransaction transaction) {
48     this(metadata, metadata.getAllocator().allocPage(transaction), keys, children, transaction);
49 }
50
51 /**
52 * Construct an inner node that is persisted to page `pageNum` allocated by
53 * metadata.getAllocator().
54 */
55 private InnerNode(BPlusTreeMetadata metadata, int pageNum, List<DataBox> keys,
56                  List<Integer> children, BaseTransaction transaction) {
57     assert(keys.size() <= 2 * metadata.getOrder());
58     assert(keys.size() + 1 == children.size());
59
60     this.metadata = metadata;
61     this.page = metadata.getAllocator().fetchPage(transaction, pageNum);
62     this.keys = keys;
63     this.children = children;
64     sync(transaction);
65 }
66
67 // Core API //////////////////////////////////////
68 // See BPlusNode.get.
69 @Override
70 public LeafNode get(BaseTransaction transaction, DataBox key) {
71     int index = numLessThanEqual(key, keys);
72     BPlusNode child = getChild(transaction, index);
73     return child.get(transaction, key);
74 }
75
76 // See BPlusNode.getLeftmostLeaf.
77 @Override
78 public LeafNode getLeftmostLeaf(BaseTransaction transaction) {
79     assert(children.size() > 0);
80     BPlusNode child = getChild(transaction, 0);
81     return child.getLeftmostLeaf(transaction);
82 }
83
84 // See BPlusNode.put.
85 @Override
86 public Optional<Pair<DataBox, Integer>> put(BaseTransaction transaction, DataBox key, RecordId rid)
87     throws BPlusTreeException {
88     int index = numLessThanEqual(key, keys);
89     BPlusNode child = getChild(transaction, index);
90     Optional<Pair<DataBox, Integer>> o = child.put(transaction, key, rid);
91
92     // If our child didn't split, then we don't have to do anything.
93     if (!o.isPresent()) {
94         return Optional.empty();
95     }
96
97     // If our child did split, then we have to insert (a) the new key and (b)
98     // the pointer to the newly created child. For example, we might go from an
99     // index which looks like this:
100     //
101     //   +---+---+---+---+
102     //   | a | b | c | e |
103     //   +---+---+---+---+
104     //   /   |   |   | \
105     //  0    1   2   3   5
106     //
107     // to an index which looks like this:
108     //
109     //   +---+---+---+---+

```

```

//      | a | b | c | d | e |
//      +-----+
//      /   |   |   |   \
//      0   1   2   3   4   5
//
// Note that in this example, p = (d, 4).
Pair<DataBox, Integer> p = o.get();
keys.add(index, p.getFirst());
children.add(index + 1, p.getSecond());

// If we can accommodate the new key and child pointer (i.e. we don't have
// more than 2d keys), then we're done (just don't forget to sync)!
int d = metadata.getOrder();
if (keys.size() <= 2 * d) {
    sync(transaction);
    return Optional.empty();
}

// On the other hand, if we overflow (i.e. we now have 2d + 1 keys), then
// we have to split ourselves in two. Continuing the example from above
// (with order 2), we would split ourselves into the following two inner
// nodes:
//
//      left          right
//      +---+---+    +---+---+
//      | a | b |    | d | e |
//      +---+---+    +---+---+
//      /   |   \   /   |   \
//      0   1   2 3  4   5
//
// We would then return the pair (c, left).
assert(keys.size() == 2 * d + 1);
List<DataBox> leftKeys = keys.subList(0, d);
DataBox middleKey = keys.get(d);
List<DataBox> rightKeys = keys.subList(d + 1, 2 * d + 1);
List<Integer> leftChildren = children.subList(0, d + 1);
List<Integer> rightChildren = children.subList(d + 1, 2 * d + 2);

// Create right node.
InnerNode n = new InnerNode(metadata, rightKeys, rightChildren, transaction);

// Update left node.
this.keys = leftKeys;
this.children = leftChildren;
sync(transaction);

return Optional.of(new Pair<>(middleKey, n.getPage().getPageNum()));
}

// See BPlusNode.bulkLoad.
@Override
public Optional<Pair<DataBox, Integer>> bulkLoad(BaseTransaction transaction,
        Iterator<Pair<DataBox, RecordId>> data,
        float fillFactor)
    throws BPlusTreeException {
    int d = metadata.getOrder();
    while (data.hasNext() && keys.size() <= 2 * d) {
        BPlusNode rightChild = getChild(transaction, children.size() - 1);
        Optional<Pair<DataBox, Integer>> o = rightChild.bulkLoad(transaction, data, fillFactor);
        if (o.isPresent()) {
            Pair<DataBox, Integer> p = o.get();
            keys.add(keys.size(), p.getFirst());
            children.add(children.size(), p.getSecond());
        }
    }

    if (keys.size() <= 2 * d) {

```

```

177         sync(transaction);
178         return Optional.empty();
179     }
180
181     assert(keys.size() == 2 * d + 1);
182     List<DataBox> leftKeys = keys.subList(0, d);
183     DataBox middleKey = keys.get(d);
184     List<DataBox> rightKeys = keys.subList(d + 1, 2 * d + 1);
185     List<Integer> leftChildren = children.subList(0, d + 1);
186     List<Integer> rightChildren = children.subList(d + 1, 2 * d + 2);
187
188     // Create right node.
189     InnerNode n = new InnerNode(metadata, rightKeys, rightChildren, transaction);
190
191     // Update left node.
192     this.keys = leftKeys;
193     this.children = leftChildren;
194     sync(transaction);
195
196     return Optional.of(new Pair<>(middleKey, n.getPage().getPageNum()));
197 }
198
199 // See BPlusNode.remove.
200 @Override
201 public void remove(BaseTransaction transaction, DataBox key) {
202     int index = numLessThanEqual(key, keys);
203     BPlusNode child = getChild(transaction, index);
204     child.remove(transaction, key);
205 }
206
207 // Helpers //////////////////////////////////////
208 @Override
209 public Page getPage() {
210     return page;
211 }
212
213 private BPlusNode getChild(BaseTransaction transaction, int i) {
214     int pageNum = children.get(i);
215     return BPlusNode.fromBytes(transaction, metadata, pageNum);
216 }
217
218 private void sync(BaseTransaction transaction) {
219     Buffer b = page.getBuffer(transaction);
220     byte[] newBytes = toBytes();
221     byte[] bytes = new byte[newBytes.length];
222     b.get(bytes);
223     if (!Arrays.equals(bytes, newBytes)) {
224         page.getBuffer(transaction).put(toBytes());
225     }
226 }
227
228 // Just for testing.
229 List<DataBox> getKeys() {
230     return keys;
231 }
232
233 // Just for testing.
234 List<Integer> getChildren() {
235     return children;
236 }
237
238 /**
239  * Returns the largest number d such that the serialization of an InnerNode
240  * with 2d keys will fit on a single page of size `pageSizeInBytes`.
241  */
242 public static int maxOrder(int pageSizeInBytes, Type keySchema) {
243     // A leaf node with n entries takes up the following number of bytes:

```

```

1 //
2 // 1 + 4 + (n * keySize) + ((n + 1) * 4)
3 //
4 // where
5 //
6 // - 1 is the number of bytes used to store isleaf,
7 // - 4 is the number of bytes used to store n,
8 // - keySize is the number of bytes used to store a DataBox of type
9 //   keySchema, and
10 // - 4 is the number of bytes used to store a child pointer.
11 //
12 // Solving the following equation
13 //
14 // 5 + (n * keySize) + ((n + 1) * 4) <= pageSizeInBytes
15 //
16 // we get
17 //
18 // n = (pageSizeInBytes - 9) / (keySize + 4)
19 //
20 // The order d is half of n.
21 int keySize = keySchema.getSizeInBytes();
22 int n = (pageSizeInBytes - 9) / (keySize + 4);
23 return n / 2;
24 }
25
26 /**
27  * Given a list ys sorted in ascending order, numLessThanEqual(x, ys) returns
28  * the number of elements in ys that are less than or equal to x. For
29  * example,
30  *
31  * numLessThanEqual(0, Arrays.asList(1, 2, 3, 4, 5)) == 0
32  * numLessThanEqual(1, Arrays.asList(1, 2, 3, 4, 5)) == 1
33  * numLessThanEqual(2, Arrays.asList(1, 2, 3, 4, 5)) == 2
34  * numLessThanEqual(3, Arrays.asList(1, 2, 3, 4, 5)) == 3
35  * numLessThanEqual(4, Arrays.asList(1, 2, 3, 4, 5)) == 4
36  * numLessThanEqual(5, Arrays.asList(1, 2, 3, 4, 5)) == 5
37  * numLessThanEqual(6, Arrays.asList(1, 2, 3, 4, 5)) == 5
38  *
39  * This helper function is useful when we're navigating down a B+ tree and
40  * need to decide which child to visit. For example, imagine an index node
41  * with the following 4 keys and 5 children pointers:
42  *
43  * +---+---+---+---+
44  * | a | b | c | d |
45  * +---+---+---+---+
46  * /   |   |   | \
47  * 0   1   2   3   4
48  *
49  * If we're searching the tree for value c, then we need to visit child 3.
50  * Not coincidentally, there are also 3 values less than or equal to c (i.e.
51  * a, b, c).
52  */
53 public static <T extends Comparable<T>> int numLessThanEqual(T x, List<T> ys) {
54     int n = 0;
55     for (T y : ys) {
56         if (y.compareTo(x) <= 0) {
57             ++n;
58         } else {
59             break;
60         }
61     }
62     return n;
63 }
64
65 /** Same as numLessThanEqual but for < instead of <= */
66 public static <T extends Comparable<T>> int numLessThan(T x, List<T> ys) {
67     int n = 0;

```

```

311     for (T y : ys) {
312         if (y.compareTo(x) < 0) {
313             ++n;
314         } else {
315             break;
316         }
317     }
318     return n;
319 }
320
321 // Pretty Printing //////////////////////////////////////
322 @Override
323 public String toString() {
324     String s = "(";
325     for (int i = 0; i < keys.size(); ++i) {
326         s += children.get(i) + " " + keys.get(i) + " ";
327     }
328     s += children.get(children.size() - 1) + ")";
329     return s;
330 }
331
332 @Override
333 public String toSexp(BaseTransaction transaction) {
334     String s = "(";
335     for (int i = 0; i < keys.size(); ++i) {
336         s += getChild(transaction, i).toSexp(transaction);
337         s += " " + keys.get(i) + " ";
338     }
339     s += getChild(transaction, children.size() - 1).toSexp(transaction) + ")";
340     return s;
341 }
342
343 /**
344  * An inner node on page 0 with a single key k and two children on page 1 and
345  * 2 is turned into the following DOT fragment:
346  *
347  * node0[label = "<f0>|k|<f1>"];
348  * ... // children
349  * "node0":f0 -> "node1";
350  * "node0":f1 -> "node2";
351  */
352 @Override
353 public String toDot(BaseTransaction transaction) {
354     List<String> ss = new ArrayList<>();
355     for (int i = 0; i < keys.size(); ++i) {
356         ss.add(String.format("<f%d>", i));
357         ss.add(keys.get(i).toString());
358     }
359     ss.add(String.format("<f%d>", keys.size()));
360
361     int pageNum = getPage().getPageNum();
362     String s = String.join("|", ss);
363     String node = String.format(" node%d[label = \"%s\\\";", pageNum, s);
364
365     List<String> lines = new ArrayList<>();
366     lines.add(node);
367     for (int i = 0; i < children.size(); ++i) {
368         BPlusNode child = getChild(transaction, i);
369         int childPageNum = child.getPage().getPageNum();
370         lines.add(child.toDot(transaction));
371         lines.add(String.format(" \\\"node%d\\\":f%d -> \\\"node%d\\\";",
372                                 pageNum, i, childPageNum));
373     }
374
375     return String.join("\\n", lines);
376 }

```

```

// Serialization //////////////////////////////////////
@Override
public byte[] toBytes() {
    // When we serialize an inner node, we write:
    //
    //   a. the literal value 0 (1 byte) which indicates that this node is not
    //       a leaf node,
    //   b. the number n (4 bytes) of keys this inner node contains (which is
    //       one fewer than the number of children pointers),
    //   c. the n keys, and
    //   d. the n+1 children pointers.
    //
    // For example, the following bytes:
    //
    //   +-----+-----+-----+-----+
    //   | 00 | 00 00 00 01 | 01 | 00 00 00 03 | 00 00 00 07 |
    //   +-----+-----+-----+-----+
    //   | \ / | \ / | \ / | \ / |
    //   |  a  |  b  |  c  |  d  |
    //
    // represent an inner node with one key (i.e. 1) and two children pointers
    // (i.e. page 3 and page 7).

    // All sizes are in bytes.
    int isLeafSize = 1;
    int numKeysSize = Integer.BYTES;
    int keysSize = metadata.getKeySchema().getSizeInBytes() * keys.size();
    int childrenSize = Integer.BYTES * children.size();
    int size = isLeafSize + numKeysSize + keysSize + childrenSize;

    ByteBuffer buf = ByteBuffer.allocate(size);
    buf.put((byte) 0);
    buf.putInt(keys.size());
    for (DataBox key : keys) {
        buf.put(key.toBytes());
    }
    for (Integer child : children) {
        buf.putInt(child);
    }
    return buf.array();
}

/**
 * InnerNode.fromBytes(t, meta, p) loads a InnerNode from page p of
 * meta.getAllocator().
 */
public static InnerNode fromBytes(BaseTransaction transaction, BPlusTreeMetadata metadata,
    int pageNum) {
    Page page = metadata.getAllocator().fetchPage(transaction, pageNum);
    Buffer buf = page.getBuffer(transaction);

    assert(buf.get() == (byte) 0);

    List<DataBox> keys = new ArrayList<>();
    List<Integer> children = new ArrayList<>();
    int n = buf.getInt();
    for (int i = 0; i < n; ++i) {
        keys.add(DataBox.fromBytes(buf, metadata.getKeySchema()));
    }
    for (int i = 0; i < n + 1; ++i) {
        children.add(buf.getInt());
    }
    return new InnerNode(metadata, pageNum, keys, children, transaction);
}

// Builtins //////////////////////////////////////
@Override

```

```
446 public boolean equals(Object o) {
447     if (o == this) {
448         return true;
449     }
450     if (!(o instanceof InnerNode)) {
451         return false;
452     }
453     InnerNode n = (InnerNode) o;
454     return page.getPageNum() == n.page.getPageNum() &&
455         keys.equals(n.keys) &&
456         children.equals(n.children);
457 }
458
459 @Override
460 public int hashCode() {
461     return Objects.hash(page.getPageNum(), keys, children);
462 }
463 }
```