

Notes on B*Trees

Thomas Dinsdale-Young, Philippa Gardner and Mark Wheelhouse

May 14, 2009

1 The Model

1.1 B^{Link}-Trees

Definition: Nodes

We define three types of node. Leaf nodes, which exist at the fringe of the tree and contain pointers to data, Intermediate nodes that contain pointers to other nodes, and Nodes where we don't know, or care, which kind of node we are dealing with.

Leaf Node : **leaf**($t_{id}, v_0, s, v_{i+1}, p_{i+1}$)
Intermediate Node : **inner**($t_{id}, v_0, p_0, s, v_{i+1}, p_{i+1}$)
Node : **node**($t_{id}, v_0, s, v_{i+1}, p_{i+1}$)

The value t_{id} tells us if the node has been locked by a thread, being set to 0 if the node is unlocked, and being set to the ID value of the locking thread otherwise. The value v_0 is the lower bound on entries in the node and the value v_{i+1} is the upper bound on entries in the node. p_0 is the pointer to a subtree with values less than this node. Note that this is the null pointer in the leaf case. p_{i+1} is the pointer to the next node at the current level of the tree and may be null if this is the rightmost node at the current level. s is a list of value-pointer pairs that allow for the quick search time of B^{link}-Trees.

Definition: Value-Pointer lists

$$s = [] \mid (v, p):s$$

where v is an integer key value and p is a pointer to a node or a data entry. In a leaf node all of the pointers in this list are to data entries, whilst in an inner node all of these pointers are to other tree nodes.

1.2 Program State

The program state consists of a *working heap* h , and a *variable store* s .

Definition: Heap

The heap h is a finite partial function that maps heap addresses to node data:

$$h : \text{ADDRS} \rightarrow_{\text{fin}} \text{NODES}$$

The heap is divided up such that each thread has its own private section of heap, and there is a single area of shared heap that contains the B*Tree.

Definition: Store

The store s is a set of finite partial functions that map integer variables $\text{Var}_{\mathbb{Z}} = \{i, j, \dots\}$ to integers, address variables $\text{Var}_{\text{ADDRS}} = \{p, q, \dots\}$ to heap addresses, boolean variables $\text{Var}_{\mathbb{B}} = \{b, \dots\}$ to boolean values and node data variables $\text{Var}_{\text{NODES}} = \{A, B, \dots\}$ to node data:

$$s : (\text{Var}_{\mathbb{Z}} \rightarrow_{\text{fin}} \mathbb{Z}) \times (\text{Var}_{\text{ADDRS}} \rightarrow_{\text{fin}} \text{ADDRS} \cup \{\text{null}\}) \times (\text{Var}_{\mathbb{B}} \rightarrow_{\text{fin}} \mathbb{B}) \times (\text{Var}_{\text{NODES}} \rightarrow_{\text{fin}} \text{NODES})$$

2 Basic Commands

2.1 Heap Commands

$A := \text{get}(x)$	<i>returns node data of node x</i>
$\text{put}(A, x)$	<i>overwrites node x with node data A</i>
$\text{lock}(x)$	<i>acquires lock of node x</i>
$\text{unlock}(x)$	<i>releases lock of node x</i>
$q := \text{new}()$	<i>creates a node of size $2k$ and returns pointer to new node</i>

2.2 Store Commands

$r := \text{root}(PB)$	<i>returns the top pointer from the primeBlock info stored in PB</i>
$p := \text{next}(A, v)$	<i>returns the pointer from A to follow to find v</i>
$p := \text{lookUp}(A, v)$	<i>returns the data pointer for v, or nil</i>
$\text{removePair}(A, v)$	<i>removes pair (v, p) from A</i>
$\text{addPair}(A, v, p)$	<i>adds pair (v, p) to A in correct place</i>
$v := \text{highValue}(A)$	<i>returns high value of A</i>
$v := \text{lowValue}(A)$	<i>returns low value of A</i>
$b := \text{isIn}(A, v)$	<i>returns true if v is a value in A, false otherwise</i>
$B := \text{rearrange}(A, v, p, q)$	<i>split A into A and B with $k \leq i \leq 2k$</i>

2.3 Language Commands

$P_1; P_2$	<i>Sequencing of commands</i>
$\text{if}(B)\{P_1\}\text{else}\{P_2\}$	<i>if $B = \text{true}$ then run P_1 otherwise run P_2</i>
$\text{while}(B)\{P\}$	<i>while $B = \text{true}$ run P</i>

2.4 Command Axioms

3 Algorithms

3.1 search

The search algorithm looks through the B*Tree trying to find the data entry associated with a certain key value val . When it gets to the node that would contain that key value, i.e. the node where $v_0 \leq val \leq v_{i+1}$, the algorithm tests to see if the key value is actually present in that node. If the key value is present then the algorithm returns a pointer to the data entry associated with that key value. If the key value is not present then the algorithm simply returns the **null** pointer.

```
r := search(val)  $\triangleq$  {  
    v := val;  
    movedown;  
    moveright;  
    if(isIn(A, v)){  
        r := lookUp(A, v)  
    }  
    else{  
        r := null  
    }  
}
```

where

```
movedown  $\triangleq$  {  
    PB := get(primeBlock)  
    current := root(PB)  
    A := get(current)  
    while(type(A)  $\neq$  leaf){  
        current := next(A, v)  
        A := get(current)  
    }  
}
```

and

```
moveright  $\triangleq$  {  
    while(v > highValue(A)){  
        current := next(A, v);  
        A := get(current)  
    }  
}
```

3.2 insert

The insert algorithm takes a value-pointer pair and inserts it into the B*Tree such that the search structure of the B*Tree is preserved.

3.3 delete

The delete algorithm searches through the B*Tree looking for the data entry associated with a certain key value. Once it finds the node that would contain the key value, i.e. the node where $v_0 \leq \text{val} \leq v_{i+1}$, the algorithm tests to see if the key value is actually present in that node. If the key value is present then the algorithm removes that key-value pair from the node and returns true. If the key value is not present then the algorithm simply returns false and makes no changes to the node. This operation may leave a node with less than k entries.

```
bool := delete(val)  $\triangleq$  {  
    v := val;  
    movedown;  
    moveright;  
    if(isIn(A,v)){  
        locking???  
        removePair(A,v;  
    )    bool := true  
    }  
    else{  
        bool := false  
    }  
}
```

3.4 compress

The compress algorithm works through (a level of?) the B*Tree and identifies nodes which have less than k entries. It compresses such nodes together so that the tree is now optimal for search times.

4 Rely Guarantee

4.1 Actions

There are 5 different actions (at this stage) that other threads can perform that will modify the shared state. These actions are,

```
LOCK  
UNLOCK  
DELETE  
INSERT@LEAF  
INSERT@INNER  
NEWROOT
```

4.2 LOCK

This action allows another thread to lock a node that is currently unlocked.

$$t \in T \setminus \{t_{id}\} \wedge \text{unlocked}(x) \rightsquigarrow \text{locked}(x, t)$$



4.3 UNLOCK

This action allows another thread to unlock a node which that same thread had previously locked.

$$t \in T \setminus \{t_{id}\} \wedge \text{locked}(x, t) \rightsquigarrow \text{unlocked}(x)$$



4.4 DELETE

This action allows another thread to remove a value-pointer pair from a leaf node if it holds the lock on that node. This action may leave a node with less than k entries.

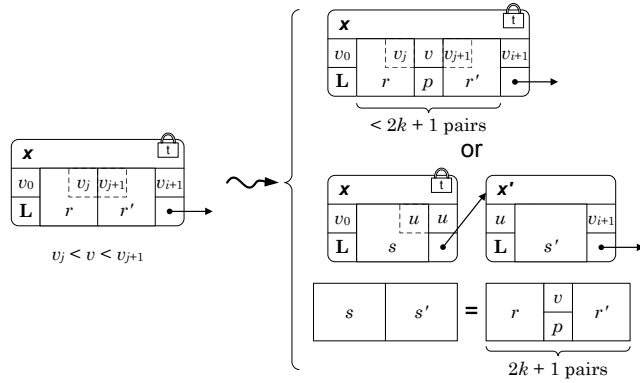
$$\begin{aligned} t \in T \setminus \{t_{id}\} \wedge \text{leaf}(t, v_0, r + (v, p) + r', v_{i+1}, p_{i+1}) \\ \rightsquigarrow \\ \text{leaf}(t, v_0, r + r', v_{i+1}, p_{i+1}) \end{aligned}$$

...diagram here...

4.5 INSERT@LEAF

This action allows another thread to add a value-pointer pair to a leaf level node if it holds the lock on that node.

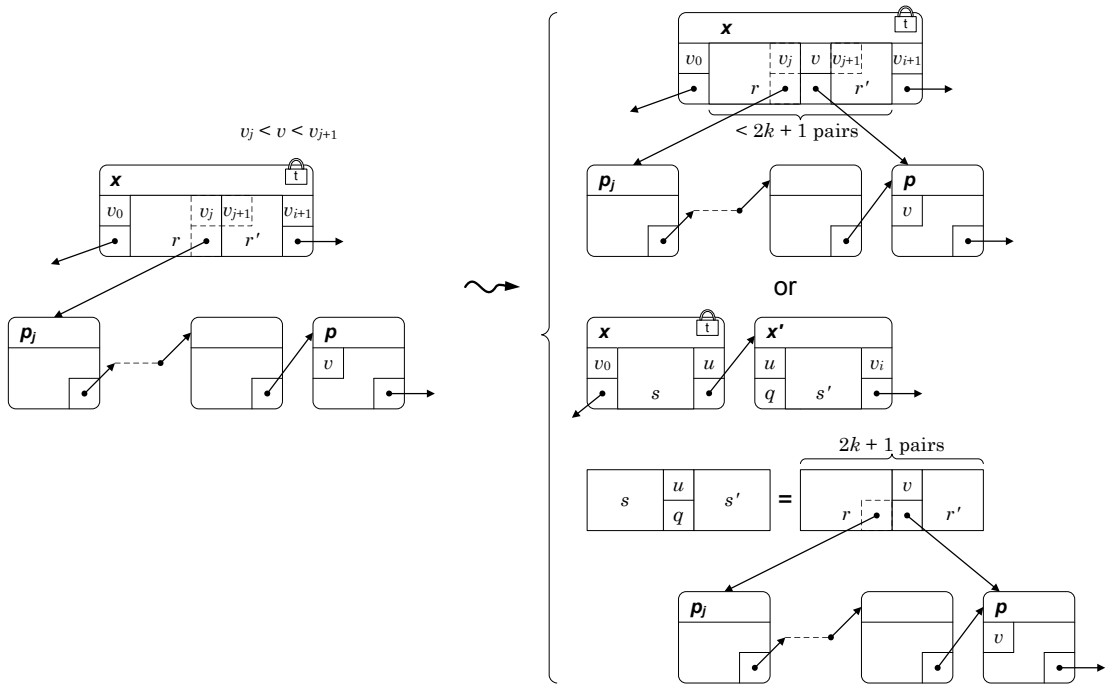
$$\begin{aligned} t \in T \setminus \{t_{id}\} \wedge \text{lastVal}(r) < v < \text{firstVal}(r') \\ \wedge x \mapsto \text{leaf}(t, v_0, r + r', v_{i+1}, p_{i+1}) \\ \rightsquigarrow \\ |r + r'| < 2k \wedge x \mapsto \text{leaf}(t, v_0, r + (v, p) + r', v_{i+1}, p_{i+1}) \\ \vee \\ |r + r'| = 2k \wedge u = \text{lastVal}(s) \wedge s + s' = r + (v, p) + r' \\ \wedge x \mapsto \text{leaf}(t, v_0, s, u, x') * x' \mapsto \text{leaf}(0, u, s', v_{i+1}, p_{i+1}) \wedge \text{created}(t, x') \end{aligned}$$



4.6 INSERT@INNER

This action allows another thread to add a value-pointer pair to an inner node if it holds the lock on that node.

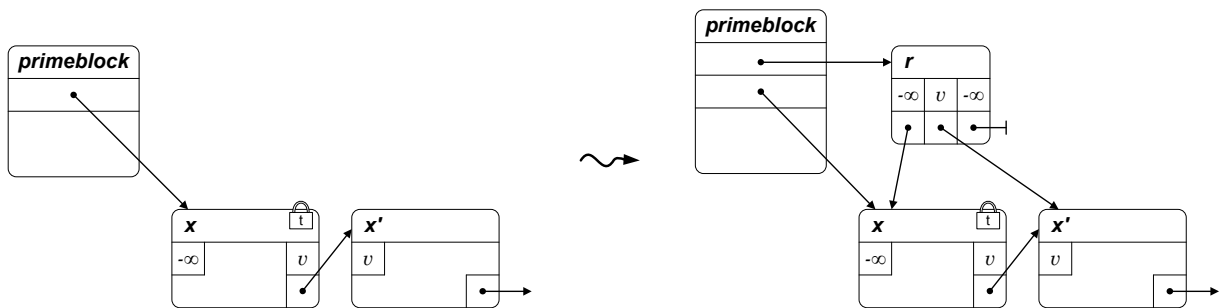
$$\begin{aligned}
& \mathbf{t} \in \mathbf{T} \setminus \{\mathbf{t}_{id}\} \wedge \text{lastVal}(\mathbf{r}) < v < \text{firstVal}(\mathbf{r}') \\
& \wedge \mathbf{x} \mapsto \text{inner}(\mathbf{t}, v_0, \mathbf{p}_0, \mathbf{r} + \mathbf{r}', v_{i+1}, \mathbf{p}_{i+1}) \\
& * \text{forest}(\mathbf{p}_j, \mathbf{p}) \wedge \mathbf{F} * \mathbf{p} \mapsto \text{node}(\neg v, \neg r, \neg r') \wedge \mathbf{P} \wedge \text{created}(\mathbf{t}, \mathbf{p}) \\
& \rightsquigarrow \\
& |\mathbf{r} + \mathbf{r}'| < 2\mathbf{k} \wedge \mathbf{x} \mapsto \text{inner}(\mathbf{t}, v_0, \mathbf{p}_0, \mathbf{r} + (v, \mathbf{p}) + \mathbf{r}', v_{i+1}, \mathbf{p}_{i+1}) * \mathbf{F} * \mathbf{P} \\
& \quad \quad \quad \downarrow \\
& |\mathbf{r} + \mathbf{r}'| = 2\mathbf{k} \wedge \mathbf{s} + (u, \mathbf{q}) + \mathbf{s}' = \mathbf{r} + (v, \mathbf{p}) + \mathbf{r}' \wedge \text{created}(\mathbf{t}, \mathbf{x}') \\
& \wedge \mathbf{x} \mapsto \text{inner}(\mathbf{t}, v_0, \mathbf{p}_0, \mathbf{s}, u, \mathbf{x}') * \mathbf{x}' \mapsto \text{inner}(0, u, \mathbf{q}, \mathbf{s}', v_{i+1}, \mathbf{p}_{i+1}) * \mathbf{F} * \mathbf{P}
\end{aligned}$$



4.7 NEWROOT

This action allows another thread to create a new root for the B*Tree if it holds the lock of the current root node.

$$\begin{aligned}
 & t \in T \setminus \{t_{id}\} \wedge \mathbf{primeBlock} \mapsto x : xs \\
 & * x \mapsto \text{node}(t, -\infty, _, v, x') \wedge X * x' \mapsto \text{node}(_, v, _, _, _) \wedge X' \wedge \text{created}(t, x') \\
 & \quad \rightsquigarrow \\
 & \mathbf{primeBlock} \mapsto r : x : xs * X * X' * r \mapsto \text{inner}(0, -\infty, x, (v, x'), \infty, \text{null})
 \end{aligned}$$



5 Correctness of Algorithms

Now that we have a notion of the interference that can be caused by other threads manipulating the B*Tree, we can look at the correctness of the algorithms provided in section 3. (So far we only have the actions that correspond to other threads performing search or insert actions, we will need to extend our action set to also consider the effects of delete and compress actions.)

We will provide specifications in a rely-guarantee style, which is of the form,

$$\mathbb{C} \models \{P\}, R, G, \{Q\}$$

This means that a program \mathbb{C} satisfies a specification if the program state satisfies P before the program initiates, the program state satisfies Q immediately after the program terminates (if it terminates), the program operates normally under the effects on the shared state specified in the Rely Set R and the program has the effects on the shared state specified by the Guarantee Set G . When we provide assertions at program points, these assertions must be stable under all of the actions allowed by the Rely Set R .

5.1 Correctness of search

To prove the correctness of search we wish to show that,

$$r := \text{search}(\text{val}) \models \{\text{BTree}\}, R', \emptyset, \left\{ \begin{array}{c} (\text{BTree} \wedge \diamond(\text{val}, p) \wedge (r = p)) \\ \vee \\ (\text{BTree} \wedge \neg \diamond(\text{val}, p) \wedge (r = \text{null})) \end{array} \right\}$$

where BTree simply states that the shared state contains a 'well-formed' B*Tree. The Rely-Set R' is the set of actions described in section 4.1 without any actions that insert (or delete) at the key value val of the search, as these would lead to a non-deterministic race for that data.

Some of the mid points from our proof (a sketch outline)

$$\exists p. \text{primeBlock} \mapsto p + \text{PB}$$

$$(\text{current} = n) \wedge n \mapsto N \diamond \text{current}$$

$$(\mathbf{A} = \mathbf{N}) \wedge \text{niceNode}(\mathbf{N}, -\infty, ?)$$

$$\text{loopInv} \wedge (\text{type}(\mathbf{A}) \neq \text{leaf})$$

$$\text{loopInv} \wedge (\text{type}(\mathbf{A}) = \text{leaf})$$

where the loop invariant is something like,

$$(\mathbf{A} = \mathbf{N}) \wedge (v' = \text{lowValue}(\mathbf{A})) \wedge \text{niceNode}(\mathbf{N}, v', ?) \wedge (v' < v)$$