

THE SPECIFICATION OF ABSTRACT MAPPINGS AND THEIR
IMPLEMENTATION AS B^+ -TREES

Elizabeth Fielding

Technical Monograph PRG-18
September 1980

Oxford University Computing Laboratory,
Programming Research Group,
45, Banbury Road,
OXFORD, OX2 6PE

ACKNOWLEDGEMENTS

I am most grateful to Cliff Jones for suggesting and supervising this project. His unfailing interest and guidance helped to make this work both a pleasant and worthwhile experience.

I also appreciate the help given to me by Jim Kaubisch and Malcolm Harper in using the UCSD system and the Sanders printer.

CONTENTS

PAGE

INTRODUCTION	1
2 THE RIGOROUS METHOD OF SPECIFICATION AND DEVELOPMENT	3
3 REPRESENTATION OF ABSTRACT MAPPINGS AS BINARY TREES	10
3.1 Specification	10
3.2 Representation 1	10
3.2.1 Proofs of correctness	12
3.3 Representation 2	17
3.4 Realization	18
4 REPRESENTATION OF ABSTRACT MAPPINGS AS B^+ -TREES	20
4.1 Description of B-trees	20
4.2 Specification	28
4.3 Representation 1	28
4.3.1 Data Structure and Data Type Invariant	29
4.3.2 The Retrieve function	30
4.3.3 The Operations: FIND1	31
INSERT1	31
DELETE1	32
4.3.4 Proofs of correctness	34
4.3.4.1 Data Type Proofs:	
Totality of the retrieve function	34
Adequacy	35

4.3.4.2 Operation Correctness Proofs:	
FIND1	36
INSERT1	38
DELETE1	47
4.4 Representation 2	55
4.4.1 Data Structure and Data Type Invariant	55
4.4.2 The Retrieve Function	57
4.4.3 The Operations: FINDP	57
INSERTP	57
DELETEP	59
4.4.4 Proofs of correctness	61
4.4.4.1 Data Type Proofs:	
Totality of the retrieve function	61
Adequacy	61
4.4.4.2 Operation Correctness Proofs:	
FINDP	62
INSERTP	64
DELETEP	69
4.5 Realization	71
5 Conclusion	73
REFERENCES	75
APPENDIX	

1. INTRODUCTION

The purpose of this project was to use the "rigorous approach" [5] to software development to develop an implementation of a type of B-tree as a structure for storing mappings from keys to data. A B-tree is a generalised binary tree, for which there exist algorithms for inserting and deleting keys which ensure that the tree remains balanced. The particular type of B-tree specified, known as a B⁺-tree, was chosen because it can allow random access to any key as well as sequential access to keys, as all keys reside in the leaves.

Initially, the simpler problem of representing abstract mappings as binary trees was tackled (and is shown) in order to illustrate and understand the entire refinement process, from abstract specification down to corresponding program code. Only the find operation was considered.

The B⁺-tree development starts with an abstract specification of a mapping from keys to data, with operations defined for finding, inserting and deleting a key. Then, two levels of specification follow, each less abstract, which represent a mapping as a tree structure and have corresponding operations for finding, insertion and deletion, which model the operations of the initial specification. The first of these levels, Representation 1, represents a tree as a set of nested sets, with the leaves of

the tree consisting of mappings from keys to data. The second level, Representation 2, represents the tree by using lists - each non-terminal node consisting of an ordered list of keys and a list of nodes, and each terminal node again consisting of a mapping from keys to data.

Each stage of the refinement is related to the preceeding stage by a retrieve function and is shown to be correct with respect to the preceeding stage in accordance with Data Type and Operation Proof Rules for Refinement [5], as described in Chapter 2. As the structure of the data types and operations of Representation 3 are similar to those of Representation 2, the correctness proofs for Representation 2 are appealed to in the arguments of the correctness of Representation 3.

The final stage of the development is the corresponding realization - running PASCAL code which implements the data types and operations. Included in the code as an indication of its correctness are weakest pre-condition type assertions [2] [3].

2. THE RIGOROUS METHOD OF SPECIFICATION AND DEVELOPMENT.

What follows has been summarized from reference [5], which should be consulted for a complete exposition.

In the rigorous method, a specification is written as a constructive specification of a data type. Development can then proceed either by operation decomposition or by data refinement. What is described below is the terminology and notation used in *constructive specification* and development by *data refinement*.

A program is considered to be an operation (or operations) on a *state* of a particular class. An *operation* can change the values of the components (or variables) that comprise a state, but cannot alter its structure.

In order to specify a program, a class of states must be defined, and it is best to design the structure of the states by choosing a data type which matches the problem as closely as possible. Such a data type is one which probably cannot be implemented directly, and this is known as an *abstract data type* and is considered to be characterized by its operations. The notation used in defining state descriptions is *abstract syntax* (for a description see [5], Ch. 14), and set, list and mapping notations are also used for describing abstract data types. An example of a *state description* is:

Studec :: N:Student-name-set Y:Student-name-set

which defines a class of states

$studc = \{ \langle n, y \rangle \mid n, y \in \text{Student-name-set} \}$

This means that any object in the class of states *studc* has two variables, with names *N* and *Y*, and in any particular state these will each be a set of student names.

A data type could be defined *implicitly*, by using axioms to relate its operations to each other, but the *constructive* approach, which is what has been used, specifies what the effects of the operations will be in terms of the underlying state.

Operations are specified by using predicates; pre- and post-conditions, as this produces shorter specifications which embody the properties required without specifying how they are to be achieved.

An operation is specified by 4 clauses, in the following format:

1) States: *s*

This specifies the name of the class of states for the operation.

2) Type: $a_1 \dots a_n \rightarrow r_1 \dots r_m$

This specifies the types of any arguments accepted and results produced.

3) pre-OP: State $a_1 \dots a_n \rightarrow \text{Bool}$

This is a predicate of a state which specifies over what subset of the class of states the operation should work.

4) post-OP: State $a_1 \dots a_n$ State $r_1 \dots r_m \rightarrow \text{Bool}$

This is a predicate of two states and defines the required relationship between the initial and final state.

An example of implicit operation definition:

(Using the state *studc*, the operation *RES* tests whether a student

RES

States: Studc

Type: Student-name \rightarrow Bool

pre-RES($\langle n, y \rangle, nm$) $\equiv nm \in (n \cup y)$

post-RES($\langle n, y \rangle, nm, \langle n', y' \rangle, b$) $\equiv n' = n$ and $y' = y$ and $b \Leftrightarrow (nm \in y)$

Now that the specification style and definitions have been defined, the method of stepwise development called *data refinement* can be described.

The initial stage is a specification which is chosen to be as abstract as possible, but must capture the required properties of the problem to be specified. A good abstraction will shorten and clarify the specification. A *data type invariant* is a predicate which must be true of all states which can be created by an operation in a specification. (Such states are said to be *valid*). The choice of data type is made so as to minimize any data type invariant required. An example of a data type invariant which might be required for states of class Studc, could state that a student name could not be a member of both the N and Y sets simultaneously;

$$\text{inv-Studc}(\langle n, y \rangle) \equiv n \cap y = \{ \}$$

The need to record a data type invariant arises because, although it may be evident from the specification and reality, it will be required explicitly in later development correctness proofs and it will also prevent errors in future alterations to the

specification. After the operations have been specified, each operation must be shown to preserve any data type invariant which might exist. The rule for *preservation of validity* is:

$$(\forall s)(\text{inv}(s) \text{ and pre-OP}(s, \text{args}) \text{ and post-OP}(s, \text{args}, s', \text{res}) \Rightarrow \text{inv}(s')) \quad (A)$$

Refinement is the term given to creating a less abstract *realization* for an abstract data type, which uses a more concrete data type, known as a *representation*, and has new operations which *model* those of the specification. Refinement is also concerned with relating a realization to its specification and *proof of its correctness*. The operations of the representation should be proved to preserve the validity according to rule (A) for the invariant and states of the representation.

Refinement is an iterative process, consisting of a series of developments, each of which is successively more concrete and is proved to be correct with respect to the stage preceeding. Should earlier invariants prove inadequate, backtracking may be required.

A *retrieve function* relates a representation to its abstraction and is the basis for data refinement proofs. Objects of a representation may contain more information than those of the abstraction and so a retrieve function operates on a state of the representation and retrieves the necessary information for the corresponding state in the abstraction.

An example of a representation for Studc might be based on a

class of states containing lists of student names:

StudCl :: NL; Student-name-list BL; Bool-list

and assuming a data type invariant:

$\text{inv-StudCl}(\langle \text{nl}, \text{bl} \rangle) \equiv \text{len nl} = \text{len bl}$

the corresponding retrieve function might be:

$\text{retr-StudCl}: \text{StudCl} \rightarrow \text{StudC}$

$\text{retr-StudCl}(\langle \text{nl}, \text{bl} \rangle) \equiv \langle \text{nl}(i) \mid 1 \leq i \leq \text{len nl} \text{ and not bl}(i) \rangle,$
 $\quad \langle \text{bl}(i) \mid 1 \leq i \leq \text{len nl} \text{ and bl}(i) \rangle$

The first part of a *refinement proof* is proof of *data type correctness*.
Two rules exist for this.

The first of these rules is used to prove the totality of the
retrieve function over valid states of the representation. It has
the form:

(Note: a suffix of 1 indicates that the function or object is an
element of the realization)

$(\forall s1)(\text{inv1}(s1) \Rightarrow (\exists s)(s = \text{retr}(s1) \text{ and inv}(\text{retr}(s1)))) \quad (\mathcal{B})$

The second of the data type proof rules is concerned with the
concept of *adequacy*. That is,

"for each (valid) element of the abstract data type there must
exist at least one value of the representation which is mapped by

the retrieve function onto the abstract value". ([5], p183)

There may exist more than one value. The formal proof rule to show adequacy is:

$$(\forall s)(\text{inv}(s) \Rightarrow (\exists s1)(\text{inv1}(s1) \text{ and } s = \text{retr}(s1))) \quad (C)$$

The second part of a refinement proof is concerned with *operation modelling*. Two rules exist for proving that each operation of the realization models the corresponding one of the specification.

The first rule, the *domains rule*, shows that the pre-condition is sufficiently wide and has the form:

$$(\forall s1)(\text{inv1}(s1) \text{ and } \text{pre-OP}(\text{retr}(s1), \text{args}) \Rightarrow \text{pre-OP1}(s1, \text{args})) \quad (D)$$

The second rule is known as the *results rule*, and states that given any state satisfying the pre-condition of OP1, and the result state after being operated on by OP1 (i.e. a state satisfying the post-condition of OP1), this pair of states must satisfy the post-condition of OP when viewed through the retrieve function. The form of this rule is:

$$(\forall s1)(\text{inv1}(s1) \text{ and } \text{pre-OP1}(s1, \text{args}) \text{ and } \text{post-OP1}(s1, \text{args}, s1', \text{res}) \Rightarrow \text{post-OP}(\text{retr}(s1), \text{args}, \text{retr}(s1'), \text{res})) \quad (E)$$

The refinement step of providing a realization and proving it correct may be repeated until a sufficiently concrete stage is

reached. As the method is *rigorous*, rather than formal, the refinement proofs need not be done formally if an informal argument can show their truth satisfactorily. In the development of B⁺-trees which follows, the proofs of the first stage of refinement are fairly formally presented whereas those of the second stage indicate how the proof could be written (often by appeal to the structure of the earlier proofs).

The final development stage is that of operation decomposition. Examples of formal proofs by weakest pre-conditions [2] [3], were followed in this step which uses "decorating assertions" in the program code as a rigorous argument of its correctness.

3. REPRESENTATION OF ABSTRACT MAPPINGS AS BINARY TREES

3.1 Specification

Since it is intended to use a binary tree as a structure in which to store a mapping from keys to data, using a mapping as the data type in the initial specification is a good abstraction. It allows the find operation to be specified in terms of its effect without prescribing how it is to work.

SPECIFICATION

$SO = \text{Key} \rightarrow \text{Data}$

APPLYO

States: SO

Type: $\text{Key} \rightarrow \text{Data}$

pre-APPLYO(s_0, k) $\cong k \in \text{dom } s_0$

post-APPLYO(s_0, k, s_0', r) $\cong s_0' = s_0$ and $r = s_0(k)$

3.2 Representation 1

This stage of the refinement represents a mapping as a binary tree. If the tree is not empty, each node of the tree contains a key and its associated data as well as a pair of pointers (either or both of which may be null), to a left and right subtree respectively. All the keys occurring in the left subtree of a particular node will have values less than the value of the key in the node, and the values of the keys of the right subtree will be greater than the node key. This is stated in the data type invariant.

$S1 = [BT1]$
 $BT1 = S1 \text{ Key Data } S1$

$invS1 : S1 \rightarrow Bool$
 $invS1(s_1) \equiv s_1 = nil \text{ or } (let \langle lt, k, d, rt \rangle = s_1 \text{ in}$
 $(\forall lk \in xks(lt)) (lk < k) \text{ and } invS1(lt) \text{ and}$
 $(\forall rk \in xks(rt)) (rk > k) \text{ and } invS1(rt))$

where $xks : S1 \rightarrow \text{Key-set}$
 $xks(s_1) \equiv \text{if } s_1 = nil$
 $\text{then } \{ \}$
 $\text{else } let \langle lt, k, d, rt \rangle = s_1 \text{ in}$
 $\{k\} \cup \text{union } \{xks(lt), xks(rt)\}$

where $union : (X\text{-set})\text{-set} \rightarrow X\text{-set}$
 $union(ss) \equiv \{e \mid \exists s \in ss \exists e \in s\}$

$retrSO : S1 \rightarrow SO$
 $retrSO(s_1) \equiv \text{if } s_1 = nil$
 $\text{then } [\]$
 $\text{else } let \langle lt, k, d, rt \rangle = s_1 \text{ in}$
 $[k \rightarrow d] \cup \text{munion } \{retrSO(lt), retrSO(rt)\}$

where $munion : (\text{Key} \rightarrow \text{Data})\text{-set} \rightarrow (\text{Key} \rightarrow \text{Data})$
 $\text{pre-munion}(ms) \equiv (\forall m1, m2 \in ms) (\text{dom } m1 \cap \text{dom } m2 = \{ \} \text{ or } m1 = m2)$
 $\text{munion}(ms) \equiv [k \rightarrow d \mid \exists m \in ms \exists k \in \text{dom } m \text{ and } d = m(k)]$

NOTE 3.1: *pre-munion* is fulfilled from *retrSO* because of the invariant *invS1*.

APPLY1
 States: $S1$
 Type: $\text{Key} \rightarrow \text{Data}$
 $\text{pre-APPLY1}(s_1, k) \equiv k \in xks(s_1)$
 $\text{post-APPLY1}(s_1, k, s_1', d) \equiv s_1' = s_1 \text{ and } d = \text{apply1}(s_1, k)$

where $\text{apply1} : BT1 \text{ Key} \rightarrow \text{Data}$
 $\text{pre-apply1}(s_1, k) \equiv k \in xks(s_1)$
 $\text{apply1}(s_1, k) \equiv \text{let } \langle lt, k', d, rt \rangle = s_1 \text{ in}$
 $\text{if } k = k'$
 $\text{then } d$
 $\text{else if } k < k'$
 $\text{then } \text{apply1}(lt, k)$
 $\text{else } \text{apply1}(rt, k)$

NOTE 3.2: *apply1* cannot be undefined because of the pre-condition, *pre-apply1*, and the invariant, *invS1*.

Two functions which are of use in the correctness proofs which follow are:

is-pdis: (X-set)-list \rightarrow Bool

is-pdis(s1) $\equiv (\forall i, j \in \{1 \dots \text{len } s1\} \mid i \neq j) \text{is-dis}(s1(i), s1(j))$

where is-dis: X-set X-set \rightarrow Bool

is-dis(s1,s2) $\equiv s1 \cap s2 = \{ \}$

3.2.1 Proofs of the Correctness of Representation 1

The refinement proofs which follow show that Representation 1 is correct with respect to the Specification given in Section 3.1.

(A) Preservation of the invariant:

The rule to be proved is:

$\text{pre-APPLY1}(s1,k) \text{ and } \text{invS1}(s1) \text{ and } \text{post-APPLY1}(s1,k,s1',d) \Rightarrow \text{invS1}(s1')$

Since APPLY1 is an identity operation on s1, the preservation of the invariant follows immediately.

(B) Totality of the retrieve function

The rule to be proved is:

$$(1) (\forall s1 \in S1)(\text{invS1}(s1) \Rightarrow (\exists s0 \in SO)(s0 = \text{retrSO}(s1)))$$

Proof: By structural induction on s1

BASIS

$$(2) s1 = nil$$

$$(3) \text{invS1}(s1)$$

$$(4) s0 = [] = \text{retrSO}(s1) \text{ which concludes the basis.}$$

INDUCTIVE HYPOTHESIS

If $s1 = \langle lt, k, d, rt \rangle$ assume

$$(5) \text{invS1}(lt) \Rightarrow (\exists s0 \in SO)(s0 = \text{retrSO}(lt))$$

$$(6) \text{invS1}(rt) \Rightarrow (\exists s0 \in SO)(s0 = \text{retrSO}(rt))$$

It follows immediately from (5) and (6), that for
 $s1 = \langle lt, k, d, rt \rangle$

since

$$(7) \text{invS1}(s1) \Rightarrow \text{is-pdisj}(\langle \text{dom retrSO}(lt), \text{dom retrSO}(rt), k \rangle)$$

that

$$(8) \text{invS1}(s1) \Rightarrow$$

$$(\exists s0 \in SO)(s0 = [k \rightarrow d] \cup \text{union} \{ \text{retrSO}(lt), \text{retrSO}(rt) \}) \quad (\text{Cf. NOTE 3.1})$$

that is

$$(9) \text{invS1}(s1) \Rightarrow (\exists s0 \in SO)(s0 = \text{retrSO}(s1))$$

(C) Adequacy

The rule to be proved is:

$$(1) (\forall s_0 \in \text{SOX}(\emptyset) \exists s_1 \in S1(\text{invS1}(s_1) \text{ and } s_0 = \text{retrSO}(s_1)))$$

Proof: By induction on $\text{dom } s_0$

BASIS

$$(2) \text{dom } s_0 = \{ \}$$

In this case

$$(3) s_0 = [\]$$

$$(4) s_1 = \text{nil} \in S1 \text{ and } \text{invS1}(s_1)$$

$$(5) s_0 = \text{retrSO}(s_1) = [\] \text{ which concludes the basis}$$

INDUCTIVE HYPOTHESIS

Assume that if

$$\text{dom } s_0 \neq \{ \} \text{ and } k \in \text{dom } s_0 \text{ then}$$

$$(6) (\forall s_0' \in \{s \in \text{SO} \mid \text{dom } s \subseteq (\text{dom } s_0 - \{k\})\}) (\exists s_1' \in S1(\text{invS1}(s_1') \text{ and } s_0' = \text{retrSO}(s_1')))$$

It now remains to be shown that s_0 can be represented.

Let

$$\{l_s, \{k\}, r_s\} \text{ be a partition of } \text{dom } s_0 \text{ such that}$$

$$(\forall e \in l_s)(k < e) \text{ and } (\forall e \in r_s)(e > k)$$

then under the induction hypothesis (6):

$$s_0 \upharpoonright l_s \text{ can be represented by } l_t \in S1 \text{ such that}$$

$$(7) \text{invS1}(l_t) \text{ and } s_0 \upharpoonright l_s = \text{retrSO}(l_t)$$

$$s_0 \upharpoonright r_s \text{ can be represented by } r_t \in S1 \text{ such that}$$

$$(8) \text{invS1}(r_t) \text{ and } s_0 \upharpoonright r_s = \text{retrSO}(r_t)$$

From this it follows that

$$s_0 \text{ can be represented by } s_1 = \langle l_t, k, s_0(k), r_t \rangle \text{ and}$$

$$(9) (\forall l_k \in \text{Xks}(l_t))(k < l_k) \text{ and } (\forall r_k \in \text{Xks}(r_t))(k < r_k)$$

$$\text{and } \text{invS1}(l_t) \text{ and } \text{invS1}(r_t)$$

$$\text{and } s_0 = [k \rightarrow d] \cup \text{union}(\text{retrSO}(l_t), \text{retrSO}(r_t))$$

which is exactly

$$(10) \text{invS1}(s_1) \text{ and } s_0 = \text{retrSO}(s_1), \text{ and this concludes the proof.}$$

Operation Proofs

(D) Domains Rule

The rule to be proved is:

$$(1) (\forall s1 \in S1)(\text{invS1}(s1) \text{ and pre-APPLY}(\text{retrSO}(s1),k) \Rightarrow \text{pre-APPLY1}(s1,k))$$

Rewriting this using the definitions of pre-APPLY and pre-APPLY1 gives:

$$(2) (\forall s1 \in S1)(\text{invS1}(s1) \text{ and } k \in \text{dom retrSO}(s1) \Rightarrow k \in \text{xks}(s1))$$

Proof: By structural induction on s1

BASIS

$$(3) s1 = nil$$

In this case (2) becomes

$$(4) k \in \{ \} \Rightarrow k \in \{ \} \text{ which is obviously true.}$$

INDUCTIVE HYPOTHESIS

If $s1 = \langle lt, k', d, rt \rangle$ then assume that

$$(5) \text{invS1}(lt) \text{ and } k \in \text{dom retrSO}(lt) \Rightarrow k \in \text{xks}(lt) \text{ and}$$

$$(6) \text{invS1}(rt) \text{ and } k \in \text{dom retrSO}(rt) \Rightarrow k \in \text{xks}(rt)$$

It now remains to be shown that (2) is true of

$$s1 = \langle lt, k', d, rt \rangle$$

$$(7) \text{invS1}(s1) \Rightarrow \text{invS1}(lt) \text{ and } \text{invS1}(rt)$$

(by definition of invS1)

$$(8) \text{invS1}(s1) \Rightarrow \text{is-pdresj}(\langle k', \text{xks}(lt), \text{xks}(rt) \rangle)$$

$$(9) k \in \text{dom retrSO}(s1) =$$

$$k \in \{k'\} \cup \text{union} \{ \text{dom retrSO}(lt), \text{dom retrSO}(rt) \}$$

$$(10) k \in \text{xks}(s1) \Leftrightarrow k \in \{k'\} \cup \text{union} \{ \text{xks}(lt), \text{xks}(rt) \}$$

Under the induction hypothesis (5) and (6),

$$(11) \text{invS1}(s1) \text{ and } k \in \text{dom retrSO}(s1) \Rightarrow k \in \text{xks}(s1) \quad \text{(by (7) and (8))}$$

(E) Results Rule

The rule to be proved is:

- (1) $(\forall s1 \in S1)(\text{invS1}(s1) \text{ and pre-APPLY1}(s1,k) \text{ and post-APPLY1}(s1,k,s1',d))$
 $\Rightarrow \text{post-APPLY}(\text{retrSO}(s1),k,\text{retrSO}(s1',d))$

Expanding this gives:

- (2) $(\forall s1 \in S1)(\text{invS1}(s1) \text{ and } k \in \text{xks}(s1) \text{ and } s1'=s1 \text{ and } d=\text{applyf1}(s1,k))$
 $\Rightarrow \text{retrSO}(s1')=\text{retrSO}(s1) \text{ and } d=\text{retrSO}(s1)(k)$

It follows immediately from $s1'=s1$

that

$$\text{retrSO}(s1')=\text{retrSO}(s1)$$

so what must be proved is:

- (3) $(\forall s1 \in S1)(\text{invS1}(s1) \text{ and } k \in \text{xks}(s1) \Rightarrow \text{applyf1}(s1,k) = \text{retrSO}(s1)(k))$

Proof: By structural induction on $s1$
BASIS

- (4) $s1 = \langle \text{nil}, k', d, \text{nil} \rangle$

In this case,

- (5) $\text{invS1}(s1) \text{ and } k \in \text{xks}(s1) \Rightarrow k=k'$
hence

- (6) $\text{applyf1}(s1,k) = \text{retrSO}(s1)(k) = d$ which proves the basis.

INDUCTIVE HYPOTHESIS

If $s1 = \langle l, k', d, r \rangle$ assume

- (7) $\text{invS1}(l) \text{ and } k \in \text{xks}(l) \Rightarrow \text{applyf1}(l,k) = \text{retrSO}(l)(k) \text{ and}$
 (8) $\text{invS1}(r) \text{ and } k \in \text{xks}(r) \Rightarrow \text{applyf1}(r,k) = \text{retrSO}(r)(k)$

Now for $s1 = \langle l, k', d, r \rangle$

Case 1:

- (9) $k = k'$ and $\text{invS1}(s1) \text{ and } k \in \text{xks}(s1)$

As for the basis, $\text{applyf1}(s1,k) = \text{retrSO}(s1)(k) = d$ (by $\text{invS1}(s1)$)

Case 2:

- (10) $k < k'$ and $\text{invS1}(s1) \text{ and } k \in \text{xks}(s1)$

In this case $k \in \text{xks}(l)$ (by $\text{invS1}(s1)$)

- (11) $\text{applyf1}(s1,k) = \text{applyf1}(l,k)$ and

- (12) $\text{retrSO}(s1)(k) = \text{retrSO}(l)(k)$

under the induction hypothesis (7) and (8) these are equivalent.

Case 3:

- (13) $k > k'$ and $\text{invS1}(s1) \text{ and } k \in \text{xks}(s1)$

This case is proved similarly to Case 2, which completes the proof of the results rule.

This level of refinement models an array, and in effect maps a binary tree onto linear storage.

REPRESENTATION 2

S2 : ROOT:[Ptr] ARRAY:Ptr \rightarrow Node2

Node2 : [Ptr] Key Data [Ptr]

invS2 S2 \rightarrow Bool

invS2(s2) \equiv (let pm = pmap(ARRAY(s2)) in
 if not (union rng pm \subseteq dom pm) then false else
 keys-are-ordered(ROOT(s2), ARRAY(s2)) and has-no-loops(pm))

where keys-are-ordered: [Ptr] (Ptr \rightarrow Node2) \rightarrow Bool

pre-keys-are-ordered(ptr,m) \equiv union rng pm \subseteq dom pm
 keys-are-ordered(ptr,m) \equiv ptr = nil or
 (let <lp,k,d,rp> = m(ptr) in
 (\forall lk \in xks2(lp,m))(lk < k) and
 (\forall rk \in xks2(rp,m))(rk > k) and
 keys-are-ordered(lp,m) and
 keys-are-ordered(rp,m))

where xks2 [Ptr] (Ptr \rightarrow Node2) \rightarrow Key-set

pre-xks2(ptr,m) \equiv union rng pm \subseteq dom pm
 xks2(ptr,m) \equiv if ptr = nil
 then { }
 else let <lp,k,d,rp> = m(ptr) in
 {k} U union {xks2(lp,m), xks2(rp,m)}

where pmap: (Ptr \rightarrow Node2) \rightarrow (Ptr \rightarrow Ptr-set)

pmap(m) \equiv [p \rightarrow (let <lp, , ,rp> = m(p) in
 {lp,rp} - { nil }) | p \in dom m]

where has-no-loops: (Ptr \rightarrow Ptr-set) \rightarrow Bool

pre-has-no-loops(pm) \equiv union rng pm \subseteq dom pm
 has-no-loops(pm) \equiv (\forall p \in dom pm)
 (p \notin union {reachable(ptr) | ptr \in pm(p)})

where reachable: Ptr (Ptr \rightarrow Ptr-set) \rightarrow Ptr-set

reachable(p,pm) \equiv build-set(p,pm,{ })

where build-set: Ptr (Ptr \rightarrow Ptr-set) Ptr-set \rightarrow Ptr-set

build-set(ptr,pm,ps) \equiv if pm(ptr) = { } or ptr \in ps
 then ps else
 union [build-set(lp,pm,{ptr} U ps) | p \in pm(ptr)]

```

retrS1: S2 → S1
retrS1(s2) ≡ retrS1(ROOT(a2), ARRAY(s2))

where retrnS1: [Ptr] (Ptr → Node2) → [Bt1]
      retrnS1(ptr,m) ≡ if ptr = nil
                    then nil
                    else let <lp,k,d,rp> = m(ptr) in
                        <retrnS1(lp,m),k,d,retrnS1(rp,m)>

```

APPLY2

States: S2

Type: Key → Data

pre-APPLY2(s2,k) ≡ k ∈ xks2(ROOT(a2), ARRAY(a2))

post-APPLY2(s2,k,a2,d) ≡ s2' = s2 and
 d = apply2(ROOT(s2),ARRAY(a2),k)

where apply2 Ptr (Ptr → Node2) Key → Data

pre-apply2(ptr,m,k) ≡ k ∈ xks2(ptr,m)

apply2(ptr,m,k) ≡ let <lp,k',d,rp> = m(ptr) in

 if k = k'

 then d

 else if k < k'

 then apply2(lp,m,k)

 else apply2(rp,m,k)

The proofs of the correctness of Representation 2 with respect to Representation 1 are not shown here. The purpose of this example of abstract maps on binary trees is only to illustrate the data refinement method which is used in specifying B⁺-trees, in which more complicated data refinement proofs are done.

3.4 Realization

This gives the PASCAL code, including assertions, corresponding to APPLY2. The implementation of the ARRAY: Ptr → Node2 component of S2 is implicit in the pointer type variables of PASCAL. An invariant, invS3, corresponding to the invariant invS2, must be true of the pointers.

APPLY3

The following data definitions are required for the procedure
apply3:

```
TYPE data_type = ... ;
   ptr = ^node;
   node = RECORD
       key:integer;
       data:data_type;
       lptr,rptr:ptr
   END;
VAR root:ptr;
    k:integer;
    res:data_type;
```

The procedure, apply3, is:

```
PROCEDURE apply3 (root:ptr; k:integer; res:data_type);
VAR p:ptr;
    node_key:integer;
BEGIN (APPLY3) (invS3(s3) and k <= xks3(root))
    p:=root;
    if p = nil
    then terminate-error
    else node_key:=p^.key;
    while node_key <> k do
    begin (node_key = p^.key and invS3(s3) and
        k <= xks3(p) and node_key # k)
        if k < node_key
        then p:=p^.lptr
        else p:=p^.rptr;
        if p = nil
        then terminate-error
        else node_key:=p^.key
    end; (node_key = p^.key and invS3(s3) and
        k <= xks3(p) and node_key = k)
    res:=p^.data
END; (APPLY3) (res = applyf2(retrS2(root),k))
```

Termination follows from depth3(p) decreasing on each iteration of
the loop (cf. has-no-loops).

4. REPRESENTATION OF ABSTRACT MAPPINGS AS B+-TREES

4.1 Description of B-trees [1] [6] [7]

A B-tree is a useful structure for storing large mappings from keys to data, where the keys are unique and have some natural order. A B-tree is a generalization of a binary tree, and a B⁺-tree is a special form of B-tree in which all the keys and data reside in the leaves. A B-tree of order m has at least m and at most 2m keys at each non-leaf node other than the root, and one more pointer to a descendant node than key in each of these nodes (i.e. between m+1 and 2m+1 pointers). The leaves must contain from m to 2m keys and, in the case of a B⁺-tree, the same number of pointers to data as keys. Thus the nodes are always at least half full. A B-tree always remains balanced - all the leaves occur at the same depth. Unless the root is a leaf, it must contain at least 1 key and 2 pointers.

A B-tree can allow the following operations to be performed:

(assume k_i denotes the i th key and d_i the associated data)

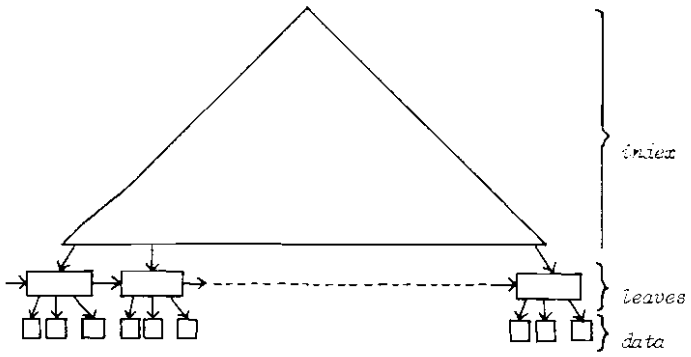
find: retrieve data d_i associated with a given key k_i

insert: add a key, k_i and its data, d_i to the mapping, provided that k_i is unique

delete: remove data d_i associated with a given key k_i

next: retrieve d_{i+1} given that d_i has just been retrieved

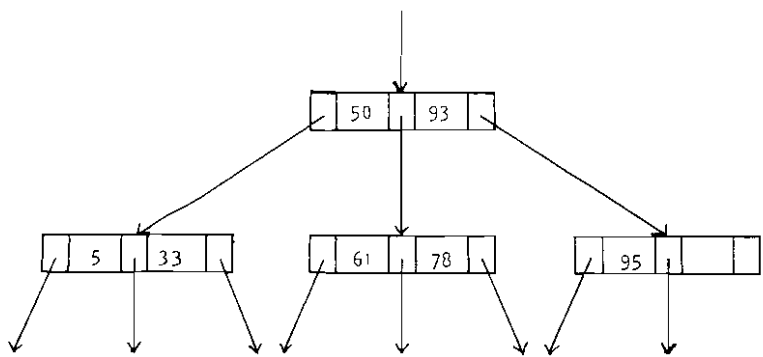
A B⁺-tree is organized into an index of non-leaf nodes, and a sequence set of leaf nodes, which may be linked sequentially from left to right as depicted, which facilitates the next operation which is laborious in an ordinary B-tree.



A B-tree

The algorithms for insertion and deletion ensure that the B-tree always remains balanced. The find, insert and delete operations are described below for a basic B-tree, with examples, and the differences in the algorithms for B⁺-trees are then given in each case.

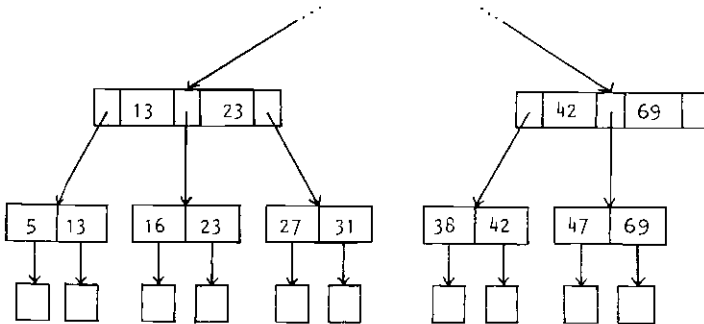
The find operation



A B-tree of order 3

Suppose that key 67 is to be found. The search starts at the root and 3 possible paths may be taken. For keys ≤ 50 the leftmost path would be taken; for keys > 50 and ≤ 93 the centre path is chosen and for keys > 93 the rightmost path is selected. This selection process is repeated at each node until an exact match is found or a leaf is reached - which denotes that the key has not been found.

For a B^+ -tree, a find operation must search all the way to a leaf, as all the keys reside in the leaves, and the key values in non-leaf nodes simply serve as separators as these nodes do not contain data.



A B^+ -tree of order 1

The insert operation

This occurs in two stages. Firstly a find operation is carried out, which must progress all the way down to the correct leaf for insertion. The insertion takes place in the leaf and the balance of the tree is restored, if necessary, by a procedure which works up from the leaf to the root. If the find stops at a leaf that is not full, the new key and data are simply inserted. If however, the leaf is full (i.e. it contains $2m$ keys) it must be *split* into two nodes with the smallest m keys and the associated data in one node,

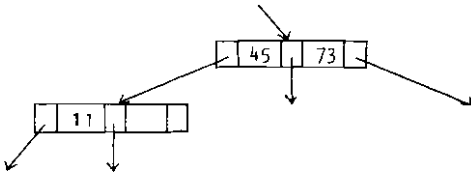
the largest m keys and data in a second node, and the middle key is inserted into the keylist of the parent node to become a separator. If the parent node is not full, the key can be added and the insertion process completed. If the parent node is full, it must be split in a similar manner. If the splitting process propagates all the way to the root, and it also has to be split, then the tree increases one level in height - it grows from the root.

In the case of a B⁺-tree, the insertion algorithm is similar, with the only difference occurring when a leaf node is split. Then, instead of the middle key being promoted to the parent node, only a copy of this key is promoted, as all the keys must reside in the leaves. Otherwise the insert operation works in the same way as for B-trees.

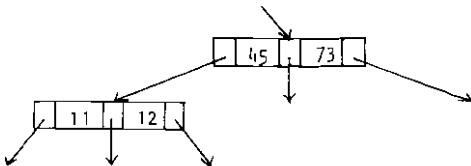
Examples

Insertion in a B-tree of order 1

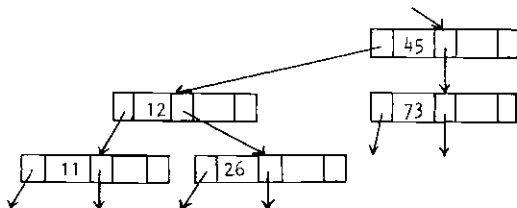
Starting with a B-tree of the form:



Insertion of the key 12 would yield:

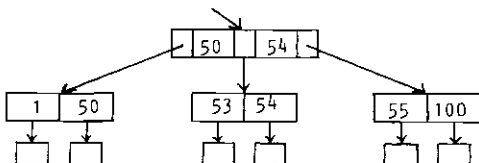


and insertion of the key 26 would produce:

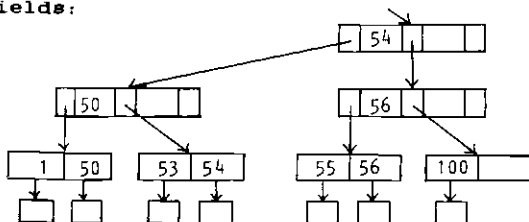


Insertion in a B⁺-tree of order 1

Insertion of the key 56 into:



yields:



The delete operation

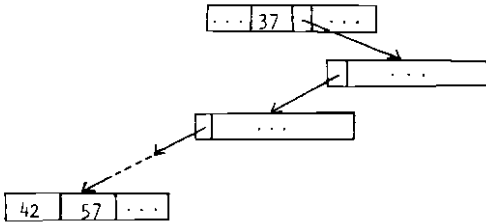
Deletion also occurs in two stages, starting with a find to locate the node containing the key to be deleted. If this key does not reside in a leaf, an adjacent key has to be found in a leaf and put into the position of the deleted key - this moves the empty position to the leaf. An adjacent key is obtained from the leftmost leaf of the right subtree of the deleted key position. If the leaf then has less than m keys, balance must be restored. If the sum of the keys of a neighbouring leaf and the leaf in question is greater than $2m$, the keys of the two nodes are evenly divided between the

nodes and the original separator key in the parent node is replaced (redistribution). If the sum of the keys is less than $2m$, the nodes are merged (the opposite of splitting) and the separator key in the parent node is pulled down and added to the combined node. If merging propagates all the way up to the root the height of the tree can decrease by one level.

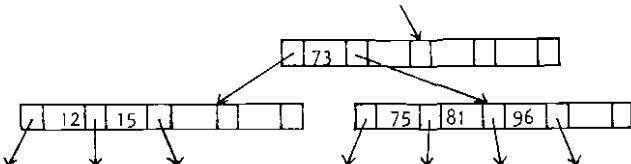
Deletion is simpler in a B⁺-tree, as non-key values may be left in the non-leaf nodes, and the key to be deleted will reside in a leaf. If a redistribution of two leaf nodes occurs, the separator key in the parent node must be overwritten with a copy of the middle key of the two nodes concerned. If a merge occurs in two leaf nodes, the separator key in the parent node is discarded. In the rest of the tree the delete operation works in exactly the same manner as in a basic B-tree.

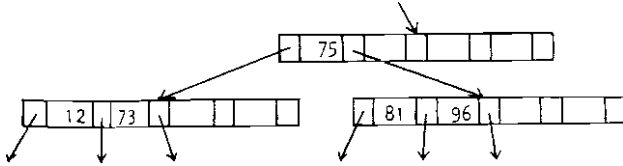
Examples
Deletion in a B-tree of order 2

The deletion of key 37 in the tree depicted below, requires that the next sequential key, 42, be found and put into the empty slot. The key to swap into position is found in the leftmost leaf of the subtree on the right of the empty slot.

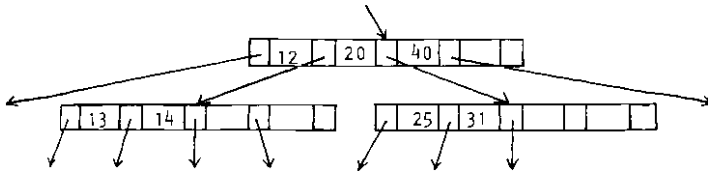


(a) Redistribution
 The deletion of key 15 in

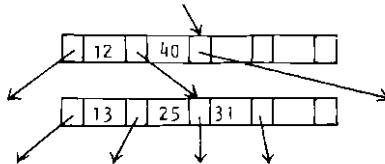




(b) Merging
 Deletion of key 14 in

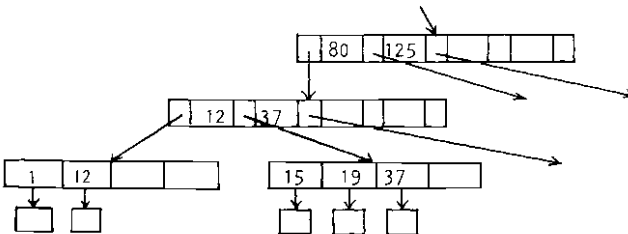


causes merging

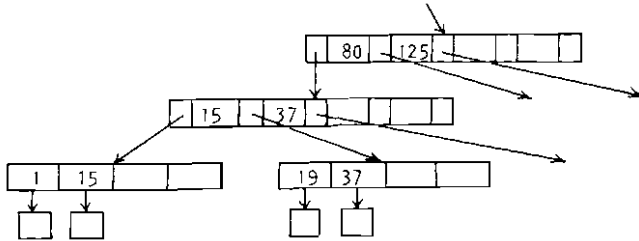


Deletion in a B⁺-tree of order 2

(a) Redistribution

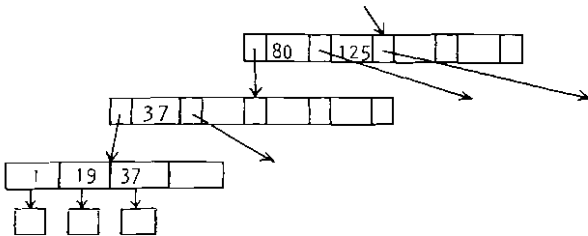


Deleting key 12 produces:



(b) Merging

Now, deleting key 15 produces:



For non-leaf nodes the deletion processes of redistribution and merging are identical to those processes acting on a B-tree, as shown in the preceding examples.

4.2 Specification

The specification uses the same abstract data type as was used for binary trees, but it has been extended by the definition of two further operations, insert and delete.

SPECIFICATION

$M = \text{Key} \rightarrow \text{Data}$

FIND

States: M

Type: $\text{Key} \rightarrow \text{Data}$

pre-FIND(m, k) $\equiv k \in \text{dom } m$

post-FIND(m, k, m', d) $\equiv m' = m \text{ and } d = m(k)$

INSERT

States: M

Type: $\text{Key Data} \rightarrow$

pre-INSERT(m, k, d) $\equiv k \notin \text{dom } m$

post-INSERT(m, k, d, m') $\equiv m' = m \cup [k \rightarrow d]$

DELETE

States: M

Type: $\text{Key} \rightarrow$

pre-DELETE(m, k) $\equiv k \in \text{dom } m$

post-DELETE(m, k, m') $\equiv m' = m \setminus \{k\}$

4.3 Representation 1

This stage uses as a representation a tree structure in the form of nested sets. This data type was chosen because it captures the essential properties of the node splitting, merging and redistributing of the B^+ -tree algorithms for insertion and deletion. However the details of the actual organization of keys in the index part of the tree does not have to be specified.

The main difficulty in the refinement proofs for this level was in proving that the insert and delete operations preserve the invariant.

4.3.1 Data Structure and Data Type Invariant

B-tree : . ORDER: Nat TREE: Node

Node = Inode | Tnode

Inode = Node-set

Tnode = Key \rightarrow Data

inv: B-tree \rightarrow Bool

inv(t) \equiv inv(ORDER(t), TREE(t))

where inv: Nat Node \rightarrow Bool

inv(m,n) \equiv common-inv(n) and size-inv(m,n)

where common-inv: Node \rightarrow Bool

common-inv(n) \equiv cases of n

n \in Inode: keysets-are-ordered(n) and balanced(n) and
(\forall sn \in n)(common-inv(sn))

n \in Tnode: true

end

where size-inv: Nat Node \rightarrow Bool

size-inv(m,n) \equiv cases of n

n \in Inode: $2 \leq \text{size}(n) \leq 2^{m+1}$ and
(\forall sn \in n)(size-inv(m,sn))

n \in Tnode: $\text{size}(n) \leq 2^m$

end

where size-inv Nat Node \rightarrow Bool

size-inv(m,n) \equiv cases of n

n \in Inode: $m+1 \leq \text{size}(n) \leq 2^{m+1}$ and
(\forall sn \in n)(size-inv(m,sn))

n \in Tnode: $m \leq \text{size}(n) \leq 2^m$

end

NOTE 4.1: The size invariant of an Inode is related to the number of descendants it can have. $m+1$ has been used so that the definitions will hold for a tree of any order ≥ 1 . The size invariant of a Tnode specifies how many keys can occur in the node.

where size: Node \rightarrow Nat

size(n) \equiv cases of n

n \in Inode: card n

n \in Tnode: card dom n

end

where **keysets-are-ordered**: Node \rightarrow Bool
keysets-are-ordered(n) \cong cases of n
 n \in Inode: ($\forall s1, s2 \in n$) ($s1 = s2$) or
 (**collect-keys**(s1) \ll **collect-keys**(s2))
 or
 (**collect-keys**(s1) \gg **collect-keys**(s2))
 n \in Tnode: true
 end

where **collect-keys**: Node \rightarrow Key-set
collect-keys(n) \cong cases of n
 n \in Inode: **union** {**collect-keys**(sn) | sn \in n}
 n \in Tnode: dom n
 end

where \ll : Nat-set Nat-set \rightarrow Bool
 $s1 \ll s2 \cong (\forall e1 \in s1 \forall e2 \in s2) (e1 < e2)$

where \gg : Nat-set Nat-set \rightarrow Bool
 $s1 \gg s2 \cong (\forall e1 \in s1 \forall e2 \in s2) (e1 > e2)$

where **union**: (X-set)-set \rightarrow X-set
union(sa) $\cong \{e \mid (\exists s \in sa) (e \in s)\}$

where **balanced**: Noda \rightarrow Bool
balanced(n) \cong card **depths**(n) = 1

where **depths**: Node \rightarrow Nat-set
depths(n) \cong cases of n
 n \in Inode: **union** {**depths**(sn) | sn \in n} ++ 1
 n \in Tnode: {1}
 end

where ++: Int-set Nat \rightarrow Int-set
 $s ++ i \cong \{e+i \mid e \in s\}$

4.3.2 The Retrieve Function

retr: B-tree \rightarrow (Key \rightarrow Data)
retr(t) \cong **retr**(TREE(t))

where **retr**: Noda \rightarrow (Key \rightarrow Data)
retr(n) \cong cases of n
 n \in Inode: **munion** {**retr**(sn) | sn \in n}
 n \in Tnode: n
 end

where **munion**: (Key \rightarrow Data)-set \rightarrow (Key \rightarrow Data)
pre-munion(ms) $\cong (\forall m1, m2 \in ms) (\text{dom } m1 \cap \text{dom } m2 = \{\} \text{ or } m1 = m2)$
munion(ms) $\cong [k \rightarrow d \mid (\exists m \in ms) (k \in \text{dom } m \text{ and } d = m(k))]$

FIND1

States: B-tree

Type: Key \rightarrow Data

pre-FIND1(t, k) $\ni k \in \text{collect-keys}(\text{TREE}(t))$

post-FIND1(t, k, t', r) $\ni t' = t$ and $r = \text{find}(k, \text{TREE}(t))$

where find: Key Node \rightarrow Data

pre-find(k, n) $\ni k \in \text{collect-keys}(n)$

find(k, n) \ni cases of n

$n \in \text{Inode}$: find($k, \text{select}(n, k)$)

$n \in \text{Tnode}$: $n(k)$

end

where select Inode Key \rightarrow Node

pre-select(n, k) $\ni k \in \text{collect-keys}(n)$

post-select(n, k, r) $\ni r \in n$ and $k \in \text{collect-keys}(r)$

INSERT1

States: B-tree

Type: Key Data \rightarrow

pre-INSERT1(t, k, d) $\ni k \in \text{collect-keys}(\text{TREE}(t))$

post-INSERT1(t, k, d, t') $\ni \text{ORDER}(t') = \text{ORDER}(t)$ and
 $\text{TREE}(t') = \text{insertn}(\text{ORDER}(t), \text{TREE}(t), k, d)$

where insertn: Nat Node Key Data \rightarrow Node

pre-insertn(m, n, k, d) $\ni k \in \text{collect-keys}(n)$

insertn(m, n, k, d) \ni if $m = \text{insertn}(m, n, k, d)$

if $\text{size}(rn) = 1$

then $\text{element}(rn)$

else rn

where insertn: Nat Node Key Data \rightarrow Inode

pre-insertn(m, n, k, d) $\ni k \in \text{collect-keys}(n)$

insertn(m, n, k, d) \ni cases of n

$n \in \text{Inode}$: if $cn = \text{select}(n, k)$

if $cns = \text{insertn}(m, cn, k, d)$

if $rn = n - \{cn\} \cup cns$

if $\text{size}(rn) \leq 2^m + 1$

then $\{rn\}$

else $\text{split}(rn)$

$n \in \text{Tnode}$: if $rn = n \cdot [k \rightarrow d]$

if $\text{size}(rn) \leq 2^m$

then $\{rn\}$

else $\text{split}(rn)$

end

where element Node-set \rightarrow Node

pre-element(ns) $\ni \text{card } ns = 1$

post-element(ns, r) $\ni r \in ns$

where selectc: Inode Key-set \rightarrow Node

post-selectc(n, ks, r) $\ni \{ \text{ins}, \text{hns} \subset n \}$ is-division($n, \text{ins}, \text{hns}, \{r\}$) and
is-ordered($\text{ins}, \text{collect-keys}(r) \cup ks, \text{hns}$)

where is-division: Inode Node-set Node-set Node-set \rightarrow Bool
 is-division(n, lns, hns, mns) \equiv is-pdij(<lns, hns, mns>) and
 (lns \cup mns \cup hns = n) (cf. section 3.2)

where is-ordered: Node-set Key-set Node-set \rightarrow Bool
 is-ordered(lns, ka, hns) \equiv (collect-keys(lns) << ka) and
 ka << collect-keys(hns)

where spliti: Inode \rightarrow Inode-set
 pre-spliti(n) \equiv size(n) \geq 2 and keysets-are-ordered(n)
 post-spliti(n, ns) \equiv union na = n and card ns = 2 and
 keysets-are-ordered(ns) and let {n1, n2} = ns in
 (size(n1) = size(n) div 2 and size(n2) = size(n) -
 size(n1))

where splitt: Tnode \rightarrow Tnode-set
 pre-splitt(n) \equiv size(n) \geq 2
 post-splitt(n, ns) \equiv union na = n and card ns = 2 and
 keysets-are-ordered(ns) and let {n1, n2} = ns in
 (size(n1) = size(n) div 2 and size(n2) = size(n) -
 size(n1))

DELETE1

States: B-tree

Type: Key \rightarrow

pre-DELETE1(l, k) \equiv k \in collect-keys(TREE(t))

post-DELETE1(l, k, t') \equiv ORDER(t') = ORDER(t) and
 TREE(t') = delete1(ORDER(t), TREE(t), k)

where delete1: Nat Node Key \rightarrow Node

pre-delete1(m, n, k) \equiv k \in collect-keys(n)

delete1(m, n, k) \equiv let rn = delete1(m, n, k)
 if m \in Inode and size(m) = 1
 then element(m)
 else m

where deleten: Nat Node Key \rightarrow Node

pre-deleten(m, n, k) \equiv k \in collect-keys(n)

deleten(m, n, k) \equiv cases of n
 n \in Inode: let cn = select(n, k)
 let rn = deleten(m, cn, k)
 if size(rn) \geq minimum-size(m, rn)
 then n - {cn} \cup {rn}
 else
 let nn = neighbour(n, cn)
 if size(rn) + size(nn) \geq
 2 * minimum-size(m, rn)
 then
 n - {cn, rn} \cup redistribute(rn, nn)
 else n - {cn, nn} \cup merge(rn, nn)
 n \in Tnode: n \ {k}
 end

where minimum-size: Nat Node \rightarrow Nat

minimum-size(m, n) \equiv cases of n
 n \in Inode: m+1
 n \in Tnode: m
 end

where neighbour: Inode Node \rightarrow Node

pre-neighbour(n,sn) $\hat{=}$ sn \in n and card n \geq 2

post-neighbour(n,sn,rn) $\hat{=}$ $\{j \mid \text{Ina}, \text{hns} \subset n\}$
(is-division(n,Ina,hns,{rn,sn}) and
is-ordered(Ina,collect-keys({rn,sn}),hns))

where redistribute: Node Node \rightarrow Inode

pre-redistribute(n1,n2) $\hat{=}$ ((n1 \in Inode and n2 \in Inode) or
(n1 \in Tnode and n2 \in Tnode)) and
size(n1) + size(n2) \geq 2 and
keys-are-ordered({n1,n2}) and
keys-are-ordered(n1) and keys-are-ordered(n2)

redistribute(n1,n2) $\hat{=}$ cases of n1,n2
n1,n2 \in Inode: split o element o merge(n1,n2)
n1,n2 \in Tnode: split o element o merge(n1,n2)
end

NOTE 4.2 This might have been better defined in terms of a post-condition for the purpose of writing the corresponding program code.

where merge: Node Node \rightarrow Inode

pre-merge(n1,n2) $\hat{=}$ (n1 \in Inode and n2 \in Inode) or
(n1 \in Tnode and n2 \in Tnode)

merge(n1,n2) $\hat{=}$ cases of n1,n2
n1,n2 \in Inode: {n1 U n2}
n1,n2 \in Tnode: {n1 + n2}
end

4.3.4 Proofs of the Correctness of Representation 1

4.3.4.1 Data Type Proofs

(B) Totality of the Retrieve Function

The rule to be proved is:

$$(1) (\forall t \in \text{B-tree})(\text{inv}(t) \Rightarrow \{\exists \text{kdm} \in (\text{Key} \rightarrow \text{Data}) \mid \text{kdm} = \text{retr}(t)\})$$

On substituting definitions, this may be written as:

$$(2) (\forall t \in \text{B-tree})(\text{inv}(\text{ORDER}(t), \text{TREE}(t)) \Rightarrow \{\exists \text{kdm} \in (\text{Key} \rightarrow \text{Data}) \mid \text{kdm} = \text{retr}(\text{TREE}(t))\})$$

Proof:

Case 1: $\text{TREE}(t) \in \text{Tnode}$

In this case, (2) becomes

$$(3) \text{size}(\text{TREE}(t)) \leq 2^{\text{ORDER}(t)} \Rightarrow \{\exists \text{kdm} \in (\text{Key} \rightarrow \text{Data}) \mid \text{kdm} = \text{TREE}(t)\}$$

which can immediately be seen to be true.

Case 2: $\text{TREE}(t) \in \text{Inode}$

In this case the proof is by structural induction on $\text{TREE}(t)$

BASIS ($\forall \text{sn} \in \text{TREE}(t) \mid \text{sn} \in \text{Tnode}$)

Equation (2) becomes

$$(4) \text{keysets-are-ordered}(\text{TREE}(t)) \text{ and } \text{balanced}(\text{TREE}(t)) \text{ and}$$

$$2 \leq \text{size}(\text{TREE}(t)) \leq 2^{\text{ORDER}(t)+1} \text{ and}$$

$$(\forall \text{sn} \in \text{TREE}(t) \mid \text{ORDER}(t) \leq \text{size}(\text{sn}) \leq 2^{\text{ORDER}(t)} \Rightarrow$$

$$\{\exists \text{kdm} \in (\text{Key} \rightarrow \text{Data}) \mid \text{kdm} = \text{union} \{ \text{sn} \mid \text{sn} \in \text{TREE}(t) \} \})$$

which can be seen to be true because *keysets-are-ordered* guarantees that the domains of the mappings to be united are disjoint.

Now suppose

($\forall \text{sn} \in \text{TREE}(t) \mid \text{sn} \in \text{Inode}$)

INDUCTIVE HYPOTHESIS

$$(5) (\forall \text{sn} \in \text{TREE}(t) \mid \text{common-inv}(\text{sn}) \text{ and } \text{size-inv}(\text{sn}) \Rightarrow$$

$$\{\exists \text{kdm} \in (\text{Key} \rightarrow \text{Data}) \mid \text{kdm} = \text{retr}(\text{sn})\})$$

In this case (2) becomes

$$(6) \text{keysets-are-ordered}(\text{TREE}(t)) \text{ and } \text{balanced}(\text{TREE}(t)) \text{ and}$$

$$(\forall \text{sn} \in \text{TREE}(t) \mid \text{common-inv}(\text{sn}) \text{ and } \text{size-inv}(\text{sn}) \text{ and}$$

$$2 \leq \text{size}(\text{TREE}(t)) \leq 2^{\text{ORDER}(t)+1} \Rightarrow$$

$$\{\exists \text{kdm} \in (\text{Key} \rightarrow \text{Data}) \mid \text{kdm} = \text{union} \{ \text{retr}(\text{sn}) \mid \text{sn} \in \text{TREE}(t) \} \})$$

which follows from the Inductive Hypothesis, because *balanced* ensures that there are no nodes with mixed Inodes and Tnodes as subnodes.

NOTE 43. The above proof could not be performed simply by induction on Node because of the special case of a tree consisting only of a single Tnode.

(C) Adequacy

The rule to be proved is:

(1) $(\forall kdm \in (\text{Key} \rightarrow \text{Data}))(\exists t \in \text{B-tree})(\text{inv}(t) \text{ and } kdm = \text{retr}(t))$

On substituting definitions, this may be written as

(2) $(\forall kdm \in (\text{Key} \rightarrow \text{Data}))(\exists t \in \text{B-tree})(\text{inv}(\text{ORDER}(t), \text{TREE}(t)) \text{ and } kdm = \text{retr}(\text{TREE}(t)))$

Proof: This is shown by an informal argument. Clearly, any mapping,

kdm such that $\text{card dom } kdm \leq 2^{\text{ORDER}(t)}$
can be represented by
 $\text{TREE}(t) = kdm$

and this satisfies
 $\text{inv}(\text{ORDER}(t), \text{TREE}(t)) \text{ and } kdm = \text{retr}(\text{TREE}(t))$

For a mapping, kdm , such that
 $\text{card dom } kdm > 2^{\text{ORDER}(t)}$

it is possible to order the domain of the mapping and then to partition kdm into a set of disjoint mappings, each of which has a domain of size between $\text{ORDER}(t)$ and $2^{\text{ORDER}(t)}$, and which form the terminal nodes of the tree. The resulting set will have cardinality > 2 and satisfies *common-inv* and *size-inv* (as $\text{ORDER}(t) > 1$). Should the set have cardinality $> 2^{\text{ORDER}(t)+1}$, it can be partitioned in such a way that each resulting set has cardinality between $\text{ORDER}(t)+1$ and $2^{\text{ORDER}(t)+1}$ and the partition has cardinality > 2 and satisfies *common-inv* and *size-inv*. This partitioning process can continue until a set results which has cardinality $\leq 2^{\text{ORDER}(t)+1}$ and what has resulted is $\text{TREE}(t)$, such that
 $\text{inv}(\text{ORDER}(t), \text{TREE}(t)) \text{ and } kdm = \text{retr}(\text{TREE}(t))$.

Adequacy can be proved formally by induction on $\text{dom } kdm$, since the proof that INSERT1 is a total operation implies that the representation of any kdm can be generated using INSERT1 .

NOTE 4.4 The special case of the root size invariant is precisely required to permit the representation of small mappings. This special case will necessitate the separation of properties several times in the proofs that follow.

4.3.4.2 Operation Proofs

Operation FIND1

(A) Preservation of the invariant

As FIND1 is an identity operation on the B-tree, it follows immediately that the invariant is preserved.

Proof that FIND1 models FIND

(D) Domains Rule

What must be shown is:

$$(1) (\forall t \in \text{B-tree}(\text{inv}(t)) \text{ and } \text{pre-FIND}(\text{retr}(t), k) \Rightarrow \text{pre-FIND1}(t, k))$$

On substituting definitions, this becomes

$$(2) (\forall t \in \text{B-tree}(\text{inv}(t)) \text{ and } k \in \text{dom } \text{retrn}(\text{TREE}(t)) \Rightarrow k \in \text{collect-keys}(\text{TREE}(t)))$$

Proof:

The first step is to prove that

$$(3) (\forall m \in \text{Nat}, n \in \text{Node}(\text{common-inv}(n)) \text{ and } k \in \text{dom } \text{retrn}(n) \Rightarrow k \in \text{collect-keys}(n))$$

This is proved by structural induction on n
BASIS

$$(4) n \in \text{Inode} \quad (2) \text{ becomes}$$

$$(5) \text{common-inv}(n) \text{ and } k \in \text{dom } n \Rightarrow k \in \text{dom } n$$

which follows immediately.

INDUCTIVE HYPOTHESIS $n \in \text{Inode}$

assume

$$(6) (\forall sn \in n(\text{common-inv}(sn)) \text{ and } k \in \text{dom } \text{retrn}(sn) \Rightarrow k \in \text{collect-keys}(sn))$$

$$(7) \text{common-inv}(n) \text{ and } k \in \text{dom } \text{retrn}(n)$$

(Hypothesis)

$$(8) (\forall sn \in n(\text{common-inv}(sn)) \text{ and}$$

(common-inv)

$$(9) (\exists sn \in n(k \in \text{dom } \text{retrn}(sn)))$$

(retrn)

$$(10) k \in \text{collect-keys}(sn) \text{ for such an } sn$$

(by (6))

Thus

$$(11) k \in \text{collect-keys}(n)$$

(collect-keys)

from which

$$(\forall t \in \text{B-tree}(\text{inv}(t)) \text{ and } k \in \text{dom } \text{retrn}(\text{TREE}(t)) \Rightarrow$$

$$k \in \text{collect-keys}(\text{TREE}(t)))$$

follows at once.

What must be shown is:

- (1) $(\forall t \in B\text{-tree}(inv(t)) \text{ and } pre\text{-FIND}(t,k) \text{ and } post\text{-FIND}(t,k,t',d) \Rightarrow$
 $post\text{-FIND}(retr(t),k,retr(t'),d))$

On substituting definitions, this becomes

- (2) $(\forall t \in B\text{-tree}(inv(t)) \text{ and } k \in collect\text{-keys}(TREE(t)) \text{ and } t'=t \text{ and}$
 $d=find(k,TREE(t)) \Rightarrow retr(t')=retr(t) \text{ and } d=retrn(TREE(t))(k))$

It follows immediately from $t'=t$ that $retr(t')=retr(t)$,

so what must be shown is

- (3) $(\forall t \in B\text{-tree}(inv(t)) \text{ and } k \in collect\text{-keys}(TREE(t)) \Rightarrow$
 $find(k,TREE(t)) = retrn(TREE(t))(k))$

Proof:

The first step is to prove that

- (3) $(\forall m \in Nat, n \in Node)(common\text{-inv}(n) \text{ and } k \in collect\text{-keys}(n) \Rightarrow$
 $find(k,n) = retrn(n)(k))$

This is proved by structural induction on n
BASIS

- (4) $n \in Inode$ (3) becomes
(5) $common\text{-inv}(n) \text{ and } k \in dom\ n \Rightarrow n(k) = n(k)$
which follows immediately.

INDUCTIVE HYPOTHESIS $n \in Inode$
assume

- (6) $(\forall sn \in n)(common\text{-inv}(sn) \text{ and } k \in collect\text{-keys}(sn) \Rightarrow$
 $find(k,sn) = retrn(sn)(k))$

- (7) $common\text{-inv}(n) \text{ and } k \in collect\text{-keys}(n)$ (Hypothesis)

- (8) $(\forall sn \in n)(common\text{-inv}(sn))$ and (common-inv)

- (9) $(\exists sn \in n)(k \in collect\text{-keys}(sn))$ (by (7) and collect-keys)

- (10) $find(k,sn) = retrn(sn)(k)$ for such a sn (by (6))

Thus

- (11) $find(k,n) = retrn(n)(k)$ (definitions of find and retrn)

from which

- $(\forall t \in B\text{-tree}(inv(t)) \text{ and } k \in collect\text{-keys}(TREE(t)) \Rightarrow$
 $find(k,TREE(t)) = retrn(TREE(t))(k))$

follows at once.

Operation INSERT1

A function which is of use in the following proofs is:

not-in: Key Node \rightarrow Bool
not-in(k,n) $\hat{=}$ k \notin collect-keys(n)

(A) Preservation of the Invariant

The rule to be proved is

(1) pre-INSERT1(t,k) and inv(t) and post-INSERT1(t,k,t') \Rightarrow inv(t')

This may be expanded to

(2) not-in(k,TREE(t)) and inv(ORDER(t),TREE(t)) \Rightarrow
inv(ORDER(t),insert(ORDER(t),TREE(t),k,d))

The proof of (2) has been done by decomposing what has to be proved into a number of lemmas and proving those which are not immediately apparent. These lemmas and their proofs then combine to provide the proof that INSERT1 preserves the invariant.

The lemmas for insertion are listed below, and are followed by the necessary proofs from which the truth of Theorem 1.12 and thus equation (2) follows:

Insertion Lemmas for Tnodes

Assuming

$m \in \text{Nat}$ and $n \in \text{Tnode}$ and $\text{not-in}(k,n)$ and $n' = \text{insertn}(m,n,k,d)$

for the following 4 lemmas:

L1.1 $\text{size-inv}(m,n) \Rightarrow n' \in \text{Tnode-set}$ and $1 \leq \text{size}(n') \leq 2$

L1.2 $\text{keysets-are-ordered}(n')$ and $\text{collect-keys}(n') = \text{collect-keys}(n) \cup \{k\}$

L1.3 $\text{balanced}(n')$

L1.4 $\text{size-inv}(m,n) \Rightarrow (\forall sn' \in n' \times \text{size-inv}(m,sn'))$

Assuming

$m \in \text{Nat}$ and $n \in \text{Tnode}$ and $\text{not-in}(k,n)$ and $\text{size-inv}(m,n)$ and $n' = \text{insertn}(m,n,k,d)$

L1.5 $\text{size-invr}(m,n)$

Insertion Lemmas for general nodes

Assuming

$m \in \text{Nat}$ and $n \in \text{Node}$ and $\text{common-inv}(n)$ and $\text{not-in}(k,n)$ and $n' = \text{insertn}(m,n,k,d)$

for the next 4 lemmas

L1.6 $\text{size-inv}(m,n) \Rightarrow n' \in \text{Node}$ and $1 \leq \text{size}(n') \leq 2$
(from L1.1, split and insertn)

L1.7 $\text{keysets-are-ordered}(n')$ and $\text{collect-keys}(n') = \text{collect-keys}(n) \cup \{k\}$

L1.8 $\text{balanced}(n')$ and $(\forall sn' \in n' \times \text{depths}(sn') = \text{depths}(n))$

L1.9 $\text{common-inv}(n')$

Assuming

$m \in \text{Nat}$ and $n \in \text{Node}$ and $\text{size-inv}(m,n)$ and $\text{not-in}(k,n)$ and $n' = \text{insertn}(m,n,k,d)$

L1.10 $(\forall sn' \in n' \times \text{size-inv}(m,sn'))$

Assuming

$m \in \text{Nat}$ and $n \in \text{Node}$ and $\text{size-invr}(m,n)$ and $\text{not-in}(k,n)$ and $n' = \text{insertn}(m,n,k,d)$

L1.11 $\text{size-invr}(m,n')$

Theorem 1.12 $\text{invr}(m,n)$ and $\text{not-in}(k,n)$ and $n' = \text{insertn}(m,n,k,d) \Rightarrow$
 $\text{invr}(m,n')$

(from L1.9 and L1.11)

Proof of Lemma 1.7

To prove:

- (1) $n \in \text{Node}$ and $\text{common-inv}(n)$ and $\text{not-in}(k,n)$ and $n = \text{insert}(m,n,k,d) \implies$
 $\text{keys-are-ordered}(n)$ and $\text{collect-keys}(n) = \text{collect-keys}(n) \cup \{k\}$

Proof: By structural induction on n

BASIS

- (2) $n \in \text{Inode}$

which follows immediately from L1.2.

INDUCTIVE HYPOTHESIS $n \in \text{Inode}$

assume

- (3) $(\forall sn \in n)(\text{common-inv}(sn)$ and $\text{not-in}(k,sn)$ and $sn = \text{insert}(m,sn,k,d)$
 $\implies \text{keys-are-ordered}(sn)$ and $\text{collect-keys}(sn) = \text{collect-keys}(sn) \cup \{k\})$

- (4) $\text{common-inv}(n)$ and $\text{not-in}(k,n)$ (Ihypothesis)

- (5) $\text{let } cn = \text{select}(n, \{k\})$

- (6) $cn \in n$ (selectc)

- (7) $\text{common-inv}(cn)$ ((3), (4), (5) and common-inv)

- (8) $\text{not-in}(k,cn)$ ((4), (5), collect-keys)

- (9) $\text{let } cns = \text{insert}(m,cn,k,d)$

- (10) $\text{keys-are-ordered}(cns)$ and $\text{collect-keys}(cns) = \text{collect-keys}(cn) \cup \{k\}$
(by (7), (8), (9), (3))

- (11) $\text{let } m = n \cdot \{cn\} \cup cns$

- (12) $\text{keys-are-ordered}(rn)$ (by (10), (11), selectc and (4))

- (13) $\text{collect-keys}(rn) = \text{collect-keys}(n) \cup \{k\}$ (by (10) and (11))

- (14) $n' = \text{if } \text{size}(rn) \leq 2^m \text{ then } \{rn\} \text{ else } \text{split}(rn)$

The lemma holds in both cases, from the definitions of split , collect-keys and keys-are-ordered .

Proof of Lemma 1.8

This is an important property of a B-tree, which is not obvious.

To prove:

- (1) $n \in \text{Node}$ and $\text{common-inv}(n)$ and $\text{not-in}(k,n)$ and $n = \text{insert}(m,n,k,d) \Rightarrow$
 $\text{balanced}(n)$ and $(\forall sn' \in n) \text{depth}(sn') = \text{depth}(n)$

Proof: By structural induction on n
BASIS

- (2) $n \in \text{Inode}$

which follows immediately from L1.3 and L1.1

INDUCTIVE HYPOTHESIS $n \in \text{Inode}$

assume

- (3) $(\forall an \in n) (\text{common-inv}(an)$ and $\text{not-in}(k,an)$ and $an = \text{insert}(m,an,k,d)$
 $\Rightarrow \text{balanced}(an)$ and $(\forall san' \in an) (\text{depth}(san') = \text{depth}(an))$)

- (4) $\text{common-inv}(n)$ and $\text{not-in}(k,n)$ (Hypothesis)

- (5) $(\forall sn \in n) (\text{common-inv}(sn)$ and $\text{not-in}(k,sn)$) ((4), common-inv and not-in)

- (6) let $cn = \text{select}(n, [k])$

- (7) $cn \in n$ (select)

- (8) $\text{common-inv}(cn)$ and $\text{not-in}(k,cn)$ ((5) and (7))

- (9) $\text{balanced}(cn)$ (by (8))

- (10) $(\forall sn \in n) (\text{depth}(sn) = \text{depth}(cn))$

$(n \in \text{Inode}$ and $\text{balanced}(n) \Rightarrow (\forall s1, s2 \in n) (\text{depth}(s1) = \text{depth}(s2))$)

- (11) $\text{depth}(n) = \text{union} \{ \text{depth}(sn) \mid sn \in n \} + 1 = \text{depth}(cn) + 1$

- (12) let $cns = \text{insert}(m, cn, k, d)$

- (13) $\text{balanced}(cns)$ and $(\forall sn' \in cns) (\text{depth}(sn') = \text{depth}(cn))$ (by (3))

- (14) let $rn = n - \{cn\} \cup cns$

- (15) $(\forall sn \in rn) (\text{depth}(sn) = \text{depth}(cn))$ (by (10) and (13))

- (16) $\text{depth}(rn) = \text{depth}(n)$ (by (10), (11) and (15))

- (17) $\text{balanced}(rn)$ (card $\text{depth}(rn) = \text{card } \text{depth}(n) = 1$ by (16))

- (18) $n' = \text{if } \text{size}(rn) < 2^m \text{ then } \{rn\} \text{ else } \text{split}(rn)$

The lemma holds in both cases, from (16), (17) and the definition of split.

Proof of Lemma 1.9

To prove:

- (1) $n \in \text{Node}$ and $\text{common-inv}(n)$ and $\text{not-in}(k,n)$ and $n' = \text{insert}(m,n,k,d) \Rightarrow \text{common-inv}(n')$

Proof: By structural induction on n

BASIS

- (2) $n \in \text{Inode}$
(3) $\text{keys-are-ordered}(n')$ (by L1.2)
(4) $\text{balanced}(n')$ (by L1.3)
(5) $n' \in \text{Node-set}$ (by L1.1)
(6) $(\forall sn' \in n')(\text{common-inv}(sn'))$
(7) $\text{common-inv}(n')$ (by (3), (4), (6))

INDUCTIVE HYPOTHESIS $n \in \text{Inode}$

assume

- (8) $(\forall sn \in n)(\text{common-inv}(sn)$ and $\text{not-in}(k,sn)$ and $sn' = \text{insert}(m,sn,k,d) \Rightarrow \text{common-inv}(sn')$)

- (9) $\text{common-inv}(n)$ and $\text{not-in}(k,n)$ (Hypothesis)
(10) $\text{keys-are-ordered}(n')$ (by L1.7)
(11) $\text{balanced}(n')$ (by L1.8)
(12) let $cn = \text{select}(n, \{k\})$
(13) $cn \in n$ (select)
(14) $\text{common-inv}(cn)$ and $\text{not-in}(k,cn)$ ((9), (13), common-inv and not-in)
(15) let $cns = \text{insert}(m,cn,k,d)$
(16) $\text{common-inv}(cns)$ (by (14), (15), (8))
(17) let $n' = n - \{cn\} \cup cns$
(18) $(\forall sn' \in n')(\text{common-inv}(sn'))$ (by (9), (13), (16) and common-inv)
(19) $\text{common-inv}(n')$ (by (10), (11) and (16))

NOTE 4.5 Induction is required for this proof because common-inv is recursive.

Proof of Lemma 1.10

To prove:

- (1) $n \in \text{Node}$ and $\text{size-inv}(m,n)$ and $\text{not-in}(k,n)$ and $n = \text{insertn}(m,n,k,d) \Rightarrow$
($\forall sn' \in n^*(\text{size-inv}(m,sn'))$)

Proof: By structural induction on n
BASIS

- (2) $n \in \text{Inode}$
which follows immediately from L1.4.

INDUCTIVE HYPOTHESIS $n \in \text{Inode}$

assume

- (3) ($\forall sn' \in n^*(\text{size-inv}(m,sn'))$ and $\text{not-in}(k,sn')$ and $sn' = \text{insertn}(m,sn',k,d)$
 $\Rightarrow (\forall ssn'' \in sn'^*(\text{size-inv}(m,ssn''))$)

- (4) $\text{size-inv}(m,n)$ and $\text{not-in}(k,n)$ (Hypothesis)

- (5) ($\forall sn' \in n^*(\text{size-inv}(m,sn'))$ and $\text{not-in}(k,sn')$) ((4), size-inv and not-in)

- (6) $\exists! cn = \text{select}(n,k)$

- (7) $cn \in n$ (select)

- (8) $\text{size-inv}(m,cn)$ and $\text{not-in}(k,cn)$ ((5) and (7))

- (9) $\exists! cns = \text{insertn}(m,cn,k,d)$

- (10) ($\forall sn' \in cns^*(\text{size-inv}(m,sn'))$) (by (3))

- (11) $1 \leq \text{size}(cns) \leq 2$ (by L1.6)

- (12) $\exists! rn = n - \{cn\} \cup cns$

- (13) $\text{size}(n) \leq \text{size}(rn) \leq \text{size}(n) + 1$

- (14) $n' = \# \text{size}(rn) \leq 2^m$ then $\{rn\}$ *is* $\text{split}(rn)$

Case (a):

- (15) $\text{size}(rn) \leq 2^m$

- (16) ($\forall sn' \in n^*(\text{size-inv}(m,sn'))$)

follows immediately from

$$\text{size}(n) \leq \text{size}(rn) \leq 2^m \text{ and } \text{size-inv}(m,n)$$

Case (b):

- (17) $\text{size}(rn) > 2^m$

- (18) $\text{size-inv}(m,n)$ and (13) $\Rightarrow \text{size}(rn) = 2^m + 1$

- (19) ($\forall sn' \in n^*(\text{size-inv}(m,sn'))$)

follows immediately from the definition of split .

Proof of Lemma 1.11

To prove:

- (1) $n \in \text{Node}$ and $\text{size-invr}(m,n)$ and $\text{not-in}(k,n)$ and $n = \text{insert}(m,n,k,d)$
 $\Rightarrow \text{size-invr}(m,n)$

Proof:

For $n \in \text{Node}$

This case follows immediately from L1.5.

For $n \in \text{inode}$

- (2) $\text{size-invr}(m,n)$ and $\text{not-in}(k,n)$ (Hypothesis)
(3) $(\forall sn \in n)(\text{size-invr}(m,sn)$ and $\text{not-in}(k,sn))$ ((2), size-invr and not-in)
(4) let $cn = \text{select}(n,k)$
(5) $cn \in n$ (selectc)
(6) $\text{size-invr}(m,cn)$ and $\text{not-in}(k,cn)$ ((3) and (5))
(7) let $cns = \text{insert}(m,cn,k,d)$
(8) $(\forall sn' \in cns)(\text{size-invr}(m,sn'))$ (by L1.10)
(9) $1 \leq \text{size}(cns) \leq 2$ (by L1.6)
(10) let $rn = n - \{cn\} \cup cns$
(11) $\text{size}(n) \leq \text{size}(rn) \leq \text{size}(n) + 1$
(12) $(\forall sm \in rn)(\text{size-invr}(m,sm))$ (by (3), (5), (6) and (10))
(13) $n' = n$ if $\text{size}(rn) \leq 2^m$ then $\{rn\}$ also split(rn)
Case (a):
 $\text{size}(rn) \leq 2^m$
 $n' = rn$ and $\text{size-invr}(m,n')$ (by (11), (12) and size-invr)
Case (b):
 $\text{size}(rn) > 2^m$
 $n' = \text{split}(rn)$ and $\text{size-invr}(m,n')$ (by (12), split)

Proof that INSERT1 models INSERT

(D) Domains Rule

What must be shown is:

$$(1) \quad \forall t \in \text{B-tree}(\text{inv}(t)) \text{ and } \text{pre-INSERT}(\text{retr}(t), k, d) \Rightarrow \\ \text{pre-INSERT1}(t, k, d)$$

On substituting definitions, this becomes

$$(2) \quad \forall t \in \text{B-tree}(\text{inv}(t)) \text{ and } k \in \text{dom } \text{retr}(t) \Rightarrow \\ k \in \text{collect-keys}(\text{TREE}(t))$$

Proof:

This proof is identical to the proof of the domains rule for operation FIND1, except that all occurrences of $k \in \dots$ are replaced by $k \in \dots$.

(E) Results Rule

The rule to be proved is:

$$\begin{aligned} e(1) \quad & (\forall t \in B\text{-tree}(\text{inv}(t)) \text{ and pre-INSERT}(t,k,d) \text{ and post-INSERT}(t,k,d,t') \\ & \Rightarrow \text{post-INSERT}(\text{retr}(t),k,d,\text{retr}(t')) \end{aligned}$$

Expanding this gives

$$\begin{aligned} (2) \quad & (\forall t \in B\text{-tree}(\text{inv}(t)) \text{ and } k \in \text{collect-keys}(\text{TREE}(t)) \text{ and } \text{ORDER}(t') = \\ & \text{ORDER}(t) \Rightarrow \\ & \text{retr}(\text{insert}(\text{ORDER}(t),\text{TREE}(t),k,d)) = \text{retr}(\text{TREE}(t)) \cup [k \rightarrow d]; \end{aligned}$$

Proof: By structural induction on TREE(t)
BASIS

(3) TREE(t) ∈ Tnode

Equation (2) becomes

$$(4) \quad \text{inv}(t) \text{ and } k \in \text{collect-keys}(\text{TREE}(t)) \text{ and } \text{ORDER}(t') = \text{ORDER}(t) \Rightarrow \\ \text{retr}(\text{ORDER}(t),\text{TREE}(t),k,d) = \text{TREE}(t) + [k \rightarrow d]$$

$$(5) \quad \text{let } rn = \text{TREE}(t) + [k \rightarrow d]$$

$$(6) \quad \text{insert}(\text{ORDER}(t),\text{TREE}(t),k,d) =$$

if size(rn) ≤ 2*ORDER(t) then {rn} else split(rn)

Case 1:

$$\text{size}(rn) \leq 2*ORDER(t)$$

$$(7) \quad \text{insert}(\text{ORDER}(t),\text{TREE}(t),k,d) = rn = \text{TREE}(t) + [k \rightarrow d]$$

$$(8) \quad \text{retr}(rn) = \text{TREE}(t) + [k \rightarrow d]$$

Case 2:

$$\text{size}(rn) > 2*ORDER(t)$$

$$(9) \quad \text{insert}(\text{ORDER}(t),\text{TREE}(t),k,d) = \text{split}(rn) = \text{split}(\text{TREE}(t) + [k \rightarrow d])$$

$$(10) \quad \text{retr}(\text{split}(rn)) = \text{munion} \{ \text{retr}(sn) \mid sn \in \text{split}(\text{TREE}(t) + [k \rightarrow d]) \} \\ = \text{TREE}(t) + [k \rightarrow d]$$

Now if TREE(t) ∈ Inode

INDUCTIVE HYPOTHESIS Assume

$$(11) \quad (\forall sn \in \text{TREE}(t)(\text{inv}(t) \text{ and } k \in \text{collect-keys}(sn) \Rightarrow \\ (\text{retr}(\text{insert}(\text{ORDER}(t),sn,k,d)) = \text{retr}(sn) + [k \rightarrow d]))$$

$$(12) \quad \text{let } cn = \text{select}(\text{TREE}(t),\{k\})$$

$$(13) \quad cn \in \text{TREE}(t)$$

(select)

$$(14) \quad \text{let } cns = \text{insert}(\text{ORDER}(t),cn,k,d)$$

$$(15) \quad \text{retr}(cns) = \text{retr}(cn) + [k \rightarrow d]$$

(by (11))

$$(16) \quad \text{let } rn = \text{TREE}(t) - \{cn\} \cup cns$$

$$(17) \quad \text{insert}(\text{ORDER}(t),\text{TREE}(t),k,d) =$$

if size(rn) ≤ 2*ORDER(t)+1 then {rn} else split(rn)

Case (a):

$$\text{size}(rn) \leq 2*ORDER(t)+1$$

$$(18) \quad \text{insert}(\text{ORDER}(t),\text{TREE}(t),k,d) = rn = \text{TREE}(t) - \{cn\} \cup cns$$

$$(19) \quad \text{retr}(rn) = \text{retr}(\text{TREE}(t)) - \text{retr}(cn) \cup \text{retr}(cns) + [k \rightarrow d] \\ = \text{retr}(\text{TREE}(t)) + [k \rightarrow d]$$

(cn ∈ TREE(t))

Case (b):

$$\text{size}(rn) > 2*ORDER(t)+1$$

$$(20) \quad \text{insert}(\text{ORDER}(t),\text{TREE}(t),k,d) = \text{split}(rn) = \text{split}(\text{TREE}(t) - \{cn\} \cup cns)$$

$$(21) \quad \text{retr}(\text{split}(rn)) = \text{munion} \{ \text{retr}(sn) \mid sn \in \text{split}(\text{TREE}(t) - \{cn\} \cup cns) \}$$

$$= \text{retr}(\text{TREE}(t)) - \text{retr}(cn) \cup \text{retr}(cns) + [k \rightarrow d]$$

$$= \text{retr}(\text{TREE}(t)) + [k \rightarrow d] \quad (\text{cn} \in \text{TREE}(t))$$

which completes the proof.

Operation DELETE1

A function which is of use in the following refinement proofs is:

is-in: Key Node \rightarrow Bool
is(k,n) \equiv k \in collect-keys(n)

(A) Preservation of the invariant

The rule to be proved is:

(1) pre-DELETE1(L,k) and inv(L) and post-DELETE1(L,k,t') \Rightarrow inv(t')

This may be expanded to:

(2) is-in(k,TREE(L)) and invr(ORDER(L),TREE(L)) \Rightarrow
invr(ORDER(L),deleter(ORDER(L),TREE(L),k))

Once again, the proof has been done by drawing up a list of Lemmas for Deletion (and proving those which are not obviously true), from which the truth of (2) can be deduced. These lemmas correspond to the Lemmas for Insertion in form.

The lemmas are listed below and are followed by selected proofs of those lemmas which are not immediately apparent and do not follow the lines of the proof of the corresponding lemma for insertion.

Deletion Lemmas for Tnodes

Assuming

$m \in \text{Nat}$ and $n \in \text{Tnode}$ and $\text{is-in}(k,n)$ and $n' = \text{deleter}(m,n,k)$

for the following 3 lemmas:

L1.13 $\text{size-inv}(m,n) \Rightarrow n' \in \text{Tnode}$ and $\text{size}(n') = \text{size}(n) - 1$

L1.14 $\text{keys-are-ordered}(n')$ and $\text{collect-keys}(n') = \text{collect-keys}(n) - \{k\}$

L1.15 $\text{balanced}(n')$

Assuming

$m \in \text{Nat}$ and $n \in \text{Tnode}$ and $\text{is-in}(k,n)$ and $\text{size-inv}(m,n)$ and
 $n' = \text{deleter}(m,n,k)$

L1.16 $\text{size-inv}(m,n')$

Deletion Lemmas for general nodes and Inodes

Assuming

$m \in \text{Nat}$ and $\text{common-inv}(n)$ and $\text{is-in}(k,n)$ and $n' = \text{deleter}(m,n,k)$

for the following 4 lemmas

L1.17 $n \in \text{Inode}$ and $\text{size-inv}(m,n) \Rightarrow n' \in \text{Inode}$ and $\text{size}(n') \leq \text{size}(n)$
(from merge, redistribute and delete)

also assuming that

$n \in \text{Node}$

for the next 3 lemmas:

L1.18 $\text{keys-are-ordered}(n')$ and $\text{collect-keys}(n') = \text{collect-keys}(n) - \{k\}$

L1.19 $\text{balanced}(n')$ and $\text{depths}(n') = \text{depths}(n)$

L1.20 $\text{common-inv}(n')$

(proof of this follows the proof of L1.9, and is done by induction using Lemmas 1.13, 1.14, 1.15, 1.18 and 1.19)

Assuming

$m \in \text{Nat}$ and $n \in \text{Inode}$ and $\text{size-inv}(m,n)$ and $\text{is-in}(k,n)$ and
 $n' = \text{deleter}(m,n,k)$

L1.21 $\forall sn' \in n' \times \text{size-inv}(m,sn')$

Assuming

$m \in \text{Nat}$ and $n \in \text{Inode}$ and $\text{size-inv}(m,n)$ and $\text{is-in}(k,n)$ and
 $n' = \text{deleter}(m,r,n)$

L1.22 $\text{size-inv}(m,n')$

(from L1.16 and L1.17)

Theorem 1.23: $\text{inv}(m,n)$ and $\text{is-in}(k,n)$ and $n' = \text{deleter}(m,n,k) \Rightarrow \text{inv}(n')$

(from L1.20 and L1.22)

Proof of Lemma 1.18

To prove:

- (1) $n \in \text{Node}$ and $\text{common-inv}(n)$ and $\text{is-in}(k,n)$ and $n' = \text{delete}(m,n,k) \Rightarrow$
 $\text{keysets-are-ordered}(n')$ and $\text{collect-keys}(n') = \text{collect-keys}(n) - \{k\}$

Proof: By structural induction on n
BASIS

- (2) $n \in \text{Tnode}$

which follows immediately from L1.14.

INDUCTIVE HYPOTHESIS $n \in \text{Node}$

assume

- (3) $(\forall n \in n)(\text{common-inv}(sn) \text{ and } \text{is-in}(k,sn) \text{ and } sn' = \text{delete}(m,sn,k) \Rightarrow$
 $\text{keysets-are-ordered}(sn') \text{ and } \text{collect-keys}(sn') = \text{collect-keys}(sn) - \{k\})$

- (4) $\text{common-inv}(n)$ and $\text{is-in}(k,n)$ (Hypothesis)
(5) $(\forall n \in n)(\text{common-inv}(sn))$ and ((4) and common-inv)
(6) $(\exists sn \in n)(\text{is-in}(sn))$ ((4) and is-in)
(7) let $cn = \text{select}(k,n)$
(8) $cn \in n$ and $\text{is-in}(k,cn)$ (select)
(9) $\text{common-inv}(cn)$ ((5) and (8))
(10) let $rn = \text{delete}(m,cn,k)$
(11) $\text{keysets-are-ordered}(rn)$ and $\text{collect-keys}(rn) = \text{collect-keys}(cn) - \{k\}$
(by (8), (9) and (3))

Case 1:

- $\text{size}(rn) \geq \text{minimum-size}(m,rn)$
(12) $n' = n - \{cn\} \cup \{rn\}$
(13) $\text{keysets-are-ordered}(n')$ and $\text{collect-keys}(n') = \text{collect-keys}(n) - \{k\}$ (by (11))

Case 2:

- $\text{size}(rn) < \text{minimum-size}(m,rn)$
(14) let $nn = \text{neighbour}(n,cn)$ (defined because $(m+1) \geq 2$)
(15) $nn \in n$ (neighbour)

Case 2 (a)

- $\text{size}(rn) + \text{size}(nn) \geq 2^*m$
(16) $\text{collect-keys}(\text{redistribute}(rn,nn)) = \text{collect-keys}(rn) \cup \text{collect-keys}(nn)$
(17) $\text{keysets-are-ordered}(\text{redistribute}(rn,nn))$ (by definition of redistribute)
(20) $n' = n - \{cn,nn\} \cup \text{redistribute}(rn,nn)$
(21) $\text{keysets-are-ordered}(n')$ (by (4), (11), (15), (16) and (17))
(22) $\text{collect-keys}(n') = \text{collect-keys}(n) - \{k\}$ (by (11), (15), (20), (16) and collect-keys)

Case 2 (b)

- $\text{size}(rn) + \text{size}(nn) < 2^*m$
(23) $\text{collect-keys}(\text{merge}(rn,nn)) = \text{collect-keys}(rn) \cup \text{collect-keys}(nn)$
(24) $n' = n - \{cn,nn\} \cup \text{merge}(rn,nn)$
(25) $\text{keysets-are-ordered}(n')$ and (by (4), (11), (23), merge, keysets-are-ordered and neighbour)
(26) $\text{collect-keys}(n') = \text{collect-keys}(n) - \{k\}$ (by (11), (23), (24) and collect-keys)

Proof of Lemma 1.19

To prove:

- (1) $n \in \text{Node}$ and $\text{common-inv}(n)$ and $\text{is-in}(k,n)$ and $n = \text{deleten}(m,n,k) \Rightarrow$
 $\text{balanced}(n)$ and $\text{depths}(n) = \text{depths}(n)$

Proof: By structural induction on n

BASIS

- (2) $n \in \text{Inode}$

which follows immediately from L1.13 and L1.15.

INDUCTIVE HYPOTHESIS $n \in \text{Inode}$

assume

- (3) $(\forall sn \in n \setminus \text{common-inv}(sn) \text{ and } \text{is-in}(k,sn) \text{ and } sn = \text{deleten}(m,sn,k) \Rightarrow$
 $\text{balanced}(sn) \text{ and } \text{depths}(sn) = \text{depths}(sn))$

- (4) $\text{common-inv}(n)$ and $\text{is-in}(k,n)$

(Hypothesis)

- (5) $(\forall sn \in n \setminus \text{common-inv}(sn))$ and

((4) and common-inv)

- (6) $(\exists sn \in n \setminus \text{is-in}(k,sn))$

((4) and is-in)

- (7) let $cn = \text{select}(k,n)$

- (8) $cn \in n$ and $\text{is-in}(k,cn)$

(select)

- (9) $\text{common-inv}(cn)$

((5) and (8))

- (10) $(\forall sn \in n \setminus \text{depths}(sn) = \text{depths}(cn))$

$(n \in \text{Inode} \text{ and } \text{balanced}(n) \Rightarrow (\forall s1, s2 \in n \setminus \text{depths}(s1) = \text{depths}(s2)))$

- (11) $\text{depths}(n) = \text{union}(\text{depths}(sn) \mid sn \in n) + 1$
 $= \text{depths}(cn) + 1$

- (12) let $m = \text{deleten}(m,cn,k)$

- (13) $\text{balanced}(rn)$ and $\text{depths}(rn) = \text{depths}(cn)$

(by (8), (9) and (3))

Case 1:

$\text{size}(n) \geq \text{minimum-size}(m,rn)$

- (14) $n' = n - \{cn\} \cup \{rn\}$

- (15) $(\forall sn' \in n' \setminus \text{depths}(sn') = \text{depths}(cn))$

(by (10) and (13))

- (16) $\text{depths}(n') = \text{depths}(n)$

(by (10), (11), (15))

- (17) $\text{balanced}(n')$

(card $\text{depths}(n) = \text{card } \text{depths}(n') = 1$ by (4))

Case 2:

$\text{size}(rn) < \text{minimum-size}(m,rn)$

- (18) let $nn = \text{neighbour}(rn,cn)$

(defined because $(m+1) \geq 2$)

- (19) $nn \in n$

(neighbour)

- (20) $\text{depths}(nn) = \text{depths}(cn) = \text{depths}(rn)$

(by (10) and (13))

Case 2(a)

$\text{size}(n) + \text{size}(nn) \geq 2 * \text{minimum-size}(m,rn)$

- (21) $\text{depths}(\text{redistribute}(rn,nn)) = \text{depths}(rn) = \text{depths}(cn)$

(by 20)

- (22) $n' = n - \{cn,nn\} \cup \text{redistribute}(rn,nn)$

- (23) $(\forall sn' \in n' \setminus \text{depths}(sn') = \text{depths}(cn))$

(by (20) and (21))

- (24) $\text{depths}(n') = \text{depths}(n)$

(by (10) and (23))

- (25) $\text{balanced}(n')$

(by card $\text{depths}(n) = 1 = \text{card } \text{depths}(n')$, (24) and (4))

Case 2(b)

$\text{size}(rn) + \text{size}(nn) < 2 * \text{minimum-size}(m,rn)$

This follows similar arguments to those of Case 2(a), and uses the property of merge:

$\text{depths}(rn) = \text{depths}(nn)$ and $\text{balanced}(rn)$ and $\text{balanced}(nn) \Rightarrow$

$\text{balanced}(\text{merge}(rn,nn))$

Proof of Lemma 1.21

To prove:

- (1) $n \in \text{inode}$ and $\text{size-inv}(m,n)$ and $\text{is-in}(k,n)$ and $n = \text{deleten}(m,n,k) \Rightarrow$
($\forall sn' \in n'(\text{size-inv}(m,sn'))$)

Proof: By structural induction on n

- (2) $\text{size-inv}(m,n)$ and $\text{is-in}(k,n)$ (Hypothesis)
(3) ($\forall sn \in n(\text{size-inv}(m,sn))$) ((2) and size-inv)
(4) ($\exists sn \in n(\text{is-in}(k,n))$) ((2) and collect-keys)

BASIS

- (5) ($\forall an \in n(\text{sn} \in \text{Inode})$)
(6) $\text{let } cn = \text{select}(k,n)$
(7) $cn \in n$ and $\text{is-in}(k,cn)$ (select)
(8) $cn \in \text{Inode}$ and $\text{size-inv}(m,cn)$ ((2), (3) and (7))
(9) $\text{let } rn = \text{deleten}(m,cn,k)$
(10) $rn = cn \setminus \{k\}$ (deleten)
(11) $\text{size}(rn) = \text{size}(cn) - 1$ (by L1.13)

Case 1:

- $\text{size}(rn) \geq \text{minimum-size}(m,rn)$
(12) $\text{size-inv}(m,rn)$ ($\text{size}(rn) < \text{size}(cn)$ and $\text{size-inv}(cn)$ by (8))
(13) $n' = n \cdot \{cn\} \cup \{rn\}$
(14) ($\forall sn' \in n'(\text{size-inv}(m,sn'))$) (by (2), (7), (12))

Case 2:

- $\text{size}(rn) < \text{minimum-size}(m,rn)$
(15) $\text{let } nn = \text{neighbour}(n,cn)$ (defined because $(m+1) \geq 2$)
(16) $nn \in n$ and $\text{size-inv}(m,nn)$ (neighbour and (3))

Case 2 (a)

- $\text{size}(rn) + \text{size}(nn) \geq 2 * \text{minimum-size}(m,rn)$
(17) $n' = n \cdot \{cn,nn\} \cup \text{redistribute}(rn,nn)$
Immediate by definition of redistribute and (8), (11) and (16).

Case 2 (b)

- $\text{size}(rn) + \text{size}(nn) < 2 * \text{minimum-size}(m,rn)$
(18) $n' = n \cdot \{cn,nn\} \cup \text{merge}(rn,nn)$
Immediate by definition of merge and (8), (11) and (16).

($\forall sn \in n(\text{sn} \in \text{Inode})$)

INDUCTIVE HYPOTHESIS

assume

- (19) ($\forall sn \in n(\text{size-inv}(m,sn)$ and $\text{is-in}(k,sn)$ and $sn = \text{deleten}(m,sn,k)$
 $\Rightarrow (\forall ssn' \in sn'(\text{size-inv}(m,ssn'))$)

- (20) $\text{let } cn = \text{select}(k,n)$
(21) $cn \in n$ and $\text{is-in}(k,cn)$ (select)
(22) ($\forall sn \in cn(\text{size-inv}(m,sn))$) (by (3), (21) and size-inv)
(23) $\text{let } rn = \text{deleten}(m,cn,k)$
(24) ($\forall srn \in rn(\text{size-inv}(m,srn))$) (by (21), (22), and (19))

Case 1:

- $\text{size}(rn) \geq \text{minimum-size}(m,rn)$
(25) $n' = n \cdot \{cn\} \cup \{rn\}$
(26) ($\forall sn' \in n'(\text{size-inv}(m,sn'))$) (by (21), (22) and (24))

Case 2:

$size(rn) < minimum-size(m, rn)$

(27) let $nn = neighbour(rn, cn)$

(defined because $(m+1) \geq 2$)

(28) $nn \in n$

(neighbour)

(29) $\forall sn \in n \setminus \{size-in(m, sn)\}$

(by (28) and (3))

Case 2(a)

$size(n) + size(nn) \geq 2 * minimum-size(m, rn)$

(30) $n' = n - \{cn, nn\} \cup redistribute(rn, nn)$

Immediate from the definition of redistribute, (22), (24) and (29).

Case 2(b)

$size(m) + size(nn) < 2 * minimum-size(m, rn)$

(31) $n' = n - \{cn, nn\} \cup merge(rn, nn)$

Immediate from the definition of merge, (22), (24) and (29).

This concludes the proof of L1.21.

NOTE 4.6: As L1.21 is concerned with Inodes, the basis cannot simply be a Tnode.

The lemma states that after DELETE1 has operated on an Inode, the size invariant is preserved for all its subnodes.

Proof that DELETE1 models DELETE

(D) Domains Rule

What must be shown is

$$(1) (\forall t \in B\text{-tree}(inv(t)) \text{ and } pre\text{-DELETE}(retr(t),k) \Rightarrow pre\text{-DELETE}(t,k))$$

On substituting definitions, this becomes

$$(2) (\forall t \in B\text{-tree}(inv(t)) \text{ and } k \in dom\ retrn(TREE(t)) \Rightarrow is\text{-in}(k, TREE(t)))$$

Proof:

This proof is identical to the proof of the domains rule for operation FIND1.

(E) Results Rule

The rule to be proved is

$$(1) (\forall t \in B\text{-tree}(inv(t)) \text{ and } pre\text{-DELETE}(t,k) \text{ and } post\text{-DELETE}(t,k,t') \Rightarrow post\text{-DELETE}(retr(t),k,retr(t')))$$

Expanding this gives

$$(2) (\forall t \in B\text{-tree}(inv(t)) \text{ and } is\text{-in}(k, TREE(t)) \text{ and } ORDER(t') = ORDER(t) \Rightarrow retrn(deleter(ORDER(t), TREE(t), k)) = retrn(TREE(t) \setminus \{k\}))$$

Proof: By structural induction on TREE(t)
BASIS

(3) TREE(t) ∈ Inode

equation (2) becomes

$$(4) inv(t) \text{ and } k \in dom\ TREE(t) \text{ and } ORDER(t') = ORDER(t) \Rightarrow retrn(deleter(ORDER(t), TREE(t), k)) = TREE(t) \setminus \{k\}$$

This follows immediately as

$$(5) deleter(ORDER(t), TREE(t), k) = deleter(ORDER(t), TREE(t), k) = TREE(t) \setminus \{k\}$$

so

$$(6) retrn(deleter(ORDER(t), TREE(t), k)) = TREE(t) \setminus \{k\}$$

If TREE(t) ∈ Inode

INDUCTIVE HYPOTHESIS Assume

$$(7) (\forall sn \in TREE(t) \setminus \{inv(t)\} \text{ and } is\text{-in}(k, sn) \Rightarrow retrn(deleter(ORDER(t), sn, k)) = retrn(sn) \setminus \{k\})$$

$$(8) let\ cn = select(TREE(t), k)$$

$$(9) cn \in TREE(t) \text{ and } is\text{-in}(k, cn)$$

(select)

$$(10) let\ rn = deleter(ORDER(t), cn, k)$$

$$(11) retrn(rn) = retrn(cn) \setminus \{k\}$$

(by (9) and (7))

Case 1 :

- $\text{size}(rn) \geq \text{minimum-size}(\text{ORDER}(t), rn)$
(12) $\text{deleten}(\text{ORDER}(t), \text{TREE}(t), k) = \text{TREE}(t) - \{cn\} \cup \{rn\}$
(13) $\text{size}(\text{deleten}(\text{ORDER}(t), \text{TREE}(t), k)) = \text{size}(\text{TREE}(t))$ (from (12))
(14) $\text{deleter}(\text{ORDER}(t), \text{TREE}(t), k) = \text{TREE}(t) - \{cn\} \cup \{rn\}$ (from (13))
(15) $\text{retrn}(\text{deleter}(\text{ORDER}(t), \text{TREE}(t), k)) =$
 $\text{retrn}(\text{TREE}(t)) - \text{retrn}(\{cn\} \cup \{rn\})$ (by (9), (11) and (14))
 $= \text{retrn}(\text{TREE}(t)) \setminus \{k\}$

Case 2 :

- $\text{size}(rn) < \text{minimum-size}(\text{ORDER}(t), rn)$
(16) $\text{let } nn = \text{neighbour}(\text{TREE}(t), cn)$ (defined because $(m+1) \geq 2$)
(17) $nn \in \text{TREE}(t)$ (neighbour)

Case 2(a)

- $\text{size}(m) + \text{size}(nn) \geq 2 * \text{minimum-size}(\text{ORDER}(t), rn)$
(18) $\text{deleten}(\text{ORDER}(t), \text{TREE}(t), k) = \text{TREE}(t) - \{cn, nn\} \cup \text{redistribute}(rn, nn)$
(19) $\text{size}(\text{deleten}(\text{ORDER}(t), \text{TREE}(t), k)) = \text{size}(\text{TREE}(t))$ (by (9), (17), (18) and redistribute)
(20) $\text{deleter}(\text{ORDER}(t), \text{TREE}(t), k) = \text{TREE}(t) - \{cn, nn\} \cup \text{redistribute}(rn, nn)$
(21) $\text{retrn}(\text{deleter}(\text{ORDER}(t), \text{TREE}(t), k)) =$
 $\text{retrn}(\text{TREE}(t)) - \text{retrn}(\{cn\} - \text{retrn}(nn) \cup \text{retrn}(\text{redistribute}(rn, nn)))$
 $= \text{retrn}(\text{TREE}(t)) - \text{retrn}(\{cn\}) - \text{retrn}(nn) \cup \text{retrn}(rn) \cup \text{retrn}(nn)$ (redistribute)
 $= \text{retrn}(\text{TREE}(t)) - \text{retrn}(\{cn\}) - \text{retrn}(nn)$ (by (11))
 $= \text{retrn}(\text{TREE}(t)) \setminus \{k\}$

Case 2(b)

- $\text{size}(rn) + \text{size}(nn) < 2 * \text{minimum-size}(\text{ORDER}(t), rn)$
(22) $\text{deleten}(\text{ORDER}(t), \text{TREE}(t), k) = \text{TREE}(t) - \{cn, nn\} \cup \text{merge}(rn, nn)$
(23) $\text{size}(\text{deleten}(\text{ORDER}(t), \text{TREE}(t), k)) = \text{size}(\text{TREE}(t)) - 1$ (by (22) and merge)
 if $\text{size}(\text{deleten}(\text{ORDER}(t), \text{TREE}(t), k)) = 1$
(24) then $\text{TREE}(t) = \{cn, nn\}$
(25) $\text{deleter}(\text{ORDER}(t), \text{TREE}(t), k) = m \cup nn$ (merge)
(26) $\text{retrn}(\text{deleter}(\text{ORDER}(t), \text{TREE}(t), k)) = \text{retrn}(rn) \cup \text{retrn}(nn)$
 $= \text{retrn}(\{cn\}) \setminus \{k\} \cup \text{retrn}(nn)$ (by (11))
 $= \text{retrn}(\text{TREE}(t)) \setminus \{k\}$ (by (24))
else
(27) if $\text{size}(\text{deleten}(\text{ORDER}(t), \text{TREE}(t), k)) > 1$
 $\text{deleter}(\text{ORDER}(t), \text{TREE}(t), k) = \text{deleten}(\text{ORDER}(t), \text{TREE}(t), k)$ (deleter)
 $= \text{TREE}(t) - \{cn, nn\} \cup \text{merge}(rn, nn)$ (by (22))

which again gives (by a similar argument to that used in

Case 2(a) :

$$\text{retrn}(\text{deleter}(\text{ORDER}(t), \text{TREE}(t), k)) = \text{retrn}(\text{TREE}(t)) \setminus \{k\}$$

4.4 Representation 2

This stage uses as a representation a tree structure which is defined by using lists, and the actual representation of the keys in the index part of the tree is specified. An Inode now consists of a list of keys, as well as a list of nodes, and the data type invariant states in what manner these lists must be ordered.

As the data type, BP-tree, is a more concrete form of the B-tree used in Representation 1, its structure closely resembles the structure of a B-tree and many of the functions correspond closely to the functions defined for Representation 1. This facilitates the refinement proofs for this level, as they have a similar form to the corresponding proofs which were done for Representation 1.

REPRESENTATION 2

441 Data Structure and Data Type Invariant

BP-tree . ORDERP. Nat TREEP. Nodep

Nodep = Inodep | Tnodep

Inodep = KEYL Key-list TREEE: Nodep-list

Tnodep = Key \rightarrow Data

invp BP-tree \rightarrow Bool

invp(t) $\hat{=}$ invrp(ORDERP(t), TREEP(t))

where invrp: Nat Nodep \rightarrow Bool

invrp(m,n) $\hat{=}$ common-invp(n) and size-invp(m,n)

where common-invp: Nodep \rightarrow Bool

common-invp(n) $\hat{=}$ cases of n

n \in Inodep: keysets-are-orderedp(n) and balancedp(n)

and ($\forall i \in [1..lenTREEE(n)]$)

(common-invp(TREEE(nX(i))))

n \in Tnodep: true

end

where size-invrp: Nat Nodep \rightarrow Bool

size-invrp(m,n) \equiv cases of n
n \in Inodep: $2 \leq \text{sizep}(n) \leq 2^m + 1$ and
($\forall i \in \{1 \dots \text{len TREEL}(n)\}$)
(size-invrp(m,TREEL(nX*i*)))
n \in Tnodep: $\text{sizep}(n) \leq 2^m$
end

where size-invp: Nat Nodep \rightarrow Bool

size-invp(m,n) \equiv cases of n
n \in Inodep: $m+1 \leq \text{sizep}(n) \leq 2^m + 1$ and
len KEYL(n) = $\text{sizep}(n) - 1$ and
($\forall i \in \{1 \dots \text{len TREEL}(n)\}$)
(size-invp(m,TREEL(nX*i*)))
n \in Tnodep: $m \leq \text{sizep}(n) \leq 2^m$
end

where sizep: Nodep \rightarrow Nat

sizep(n) \equiv cases of n
n \in Inodep: len TREEL(n)
n \in Tnodep: card dom n
end

where keysets-are-orderedp: Nodep \rightarrow Bool

keysets-are-orderedp(n) \equiv cases of n
n \in Inodep: is-ordered(KEYL(n)) and
($\forall i \in \{1 \dots \text{len KEYL}(n)\}$)
(collect-keyssp(TREEL(nX*i*)) <<=
{KEYL(nX*i*)} <<
collect-keyssp(TREEL(nX*i*+1)))
n \in Tnodep: true
end

where is-ordered: Key-list \rightarrow Bool

is-ordered(kl) \equiv ($\forall i \in \{1 \dots (\text{len kl} - 1)\}$) (kl(*i*) < kl(*i*+1))

where collect-keyssp: Nodep \rightarrow Key-set

collect-keyssp(n) \equiv cases of n
n \in Inodep: union {collect-keyssp(TREEL(nX*i*)) |
 $1 \leq i \leq \text{len TREEL}(n)$ }
n \in Tnodep: dom n
end

where <<=: Nat-set Nat-set \rightarrow Bool

s1 <<= s2 \equiv ($\forall e1 \in s1 \forall e2 \in s2 (e1 \leq e2)$)

where balancedp: Nodep \rightarrow Bool

balancedp(n) \equiv card depthsp(n) = 1

where depthsp: Nodep \rightarrow Nat-set

depthsp(n) \equiv cases of n
n \in Inodep: union {depthsp(TREEL(nX*i*)) |
 $1 \leq i \leq \text{len TREEL}(n)$ } ++ 1
n \in Tnodep: {1}
end

4.4.2 The Retrieve Function

retrp: BP-tree \rightarrow B-tree

retrp(t) $\hat{=} \langle \text{ORDERP}(t), \text{retrnp}(\text{TREEP}(t)) \rangle$

where retrnp: Nodep \rightarrow Node

retrnp(n) $\hat{=} \text{cases of } n$
 $n \in \text{Inodep: } \{ \text{retrnp}(sn) \mid sn \in \text{elems TREEL}(n) \}$
 $n \in \text{Tnodep: } n$
end

4.4.3 The Operations

FINDP

States: BP-tree

Type: Key \rightarrow Data

pre-FINDP(t,k) $\hat{=} k \in \text{collect-keysp}(\text{TREEP}(t))$

post-FINDP(t,k,t',r) $\hat{=} t' = t$ and $r = \text{findp}(k, \text{TREEP}(t))$

where findp: Key Nodep \rightarrow Data

pre-findp(k,n) $\hat{=} k \in \text{collect-keysp}(n)$
 findp(k,n) $\hat{=} \text{cases of } n$
 $n \in \text{Inodep: } \{ \text{findp}(k, \text{TREEP}(n[i])) \mid i = \text{index}(k, \text{KEYL}(n)) \}$
 $n \in \text{Tnodep: } n(k)$
end

where index: Key Key-list \rightarrow Nat

index(k,kl) $\hat{=} \text{if } kl = \langle \rangle$
then 1
else if $k \leq \text{hd } kl$
then 1
else $1 + \text{index}(k, \text{tl } kl)$

INSERTP

States: BP-tree

Type: Key Data \rightarrow

pre-INSERTP(t,k,d) $\hat{=} k \in \text{collect-keysp}(\text{TREEP}(t))$

post-INSERTP(t,k,d,t') $\hat{=} \text{ORDERP}(t') = \text{ORDERP}(t)$ and
 $\text{TREEP}(t') = \text{insertp}(\text{ORDERP}(t), \text{TREEP}(t), k, d)$

where insertp: Nat Nodep Key Data \rightarrow Nodep

pre-insertp(m,n,k,d) $\hat{=} k \in \text{collect-keysp}(n)$
 insertp(m,n,k,d) $\hat{=} \{ \text{if } \text{rn} = \text{insertnp}(m,n,k,d)$
 $\text{if } \text{sizep}(\text{rn}) = 1$
 $\text{then } \text{TREEP}(\text{rn}(1))$
 $\text{else } \text{rn}$

where insertnp: Nat Nodep Key Data \rightarrow Inodep

pre-insertnp(m,n,k,d) $\hat{=} k \in \text{collect-keysp}(n)$
 insertnp(m,n,k,d) $\hat{=} \text{cases of } n$
 $n \in \text{Inodep: } \{ \text{insertnp}(m, n[i], k, d) \mid i, k, ck, cn = \text{selectp}(n, k)$

```

let cni = insertnp(m,cn,k,d)
let sn = << >,<cn>>
let rn = replace(n,an,cni,it)
if size(rn) < 2*m+1
then << >,<rn>>
else splitp(m,rn)
n G Tnodep: let rn = n + [k->d]
if size(rn) < 2*m
then << >,<rn>>
else splittp(m,rn)
end

```

where selectp: Inodep Key \rightarrow Nat Nat Key Nodep

```

post-selectp(n,k,rit,rik,rk,rn)  $\equiv$  rit = index(k,KEYL(n))
if rit > len KEYL(n)
then rik = rit - 1
else rik = rit
rk = KEYL(n)(rik)
rn = TREEL(n)(rit)

```

where replace: Inodep Inodep Inodep Nat Nat \rightarrow Inodep

```

pre-replace(i1,i2,i3,nk,nt)  $\equiv$  is-aubnode(i1,i2) and
((KEYL(i2)=< > and nk=nt) or
(nk=position(KEYL(i2)(1),KEYL(i1))))
and nt=position(TREEL(i2)(1),TREEL(i1))
replace(i1,i2,i3,nk,nt)  $\equiv$  <alterKEYL(i1),KEYL(i2),KEYL(i3),nk),
alter(TREEL(i1),TREEL(i2),TREEL(i3),nt)>

```

where is-subnode: Inodep Inodep \rightarrow Bool

```

is-subnode(i1,i2)  $\equiv$  if KEYL(i2) = < >
then len TREEL(i2)=1 and TREEL(i2)(1)  $\in$ 
elems TREEL(i1)
else is-aublist(KEYL(i1),KEYL(i2)) and
is-sublist(TREEL(i1,i2)) and
position(KEYL(i2)(1),KEYL(i1)) =
position(TREEL(i2)(1),TREEL(i1))

```

where is-sublist: X-list X-list \rightarrow Bool

```

pre-is-sublist(l1,l2)  $\equiv$  0 < len l2 < len l1
is-sublist(l1,l2)  $\equiv$  { }  $\in$  { 1 .. len l1 } X | position(l2)(1),1)
and (V j  $\in$  { 1 .. len l2 } X | 2)(j)=1(i+j)-1)

```

where position: X X-list \rightarrow Nat

```

pre-position(x,xl)  $\equiv$  x  $\in$  elems xl
position(x,xl)  $\equiv$  if x = hd xl
then 1
else 1+position(x, tl xl)

```

where alter: X-list X-list X-list Nat \rightarrow X-list

```

pre-alter(l1,l2,l3,i)  $\equiv$  (l2=< > and 1 < i < len l1 + 1) or
(is-aublist(l1,l2) and i=position(l2)(1),1))
alter(l1,l2,l3,i)  $\equiv$  front(i-1,l1) || l3 || back(len l1-len l2-i+1,l1)

```

NOTE 4.7: The length of the preceding five functions suggests that they could be restructured. An alternative approach would have been to return only a pair of indices to positions in the key and tree lists from selectp, and then to use only these indices as arguments in the function replace.

where front: Nat0 X-list \rightarrow X-list

front(n,l) \equiv if n=0
then < >
else if n > len l
then l
else <hd l> || front(n-1,l)

where back: Nat0 X-list \rightarrow X-list

post-back(n,l,r) \equiv if n = 0 or n > len l
then r = < >
else len r = n and
($\forall i \in [1..n]$) \times r(i) = l(len l - n + i)

where splitp: Nat Inodep \rightarrow Inodep

pre-splitp(m,n) \equiv len KEYL(n) = 2^m+1 and len TREEL(n) = 2^m+2
splitp(m,n) \equiv <<KEYL(n) \times m+1>>, <<front(m,KEYL(n)), front(m+1, TREEL(n))>>,
<<back(m,KEYL(n)), back(m+1, TREEL(n))>>>

where splittp: Nat Tnodep \rightarrow Inodep

pre-splittp(m,n) \equiv sizep(n) \geq 2
post-splittp(m,n,r) \equiv let sk,gk = halve(dom n)
r = <<max(sk), <[k \rightarrow n(k) | k \in sk],
[k \rightarrow n(k) | k \in gk]>>

where halve: Key-set \rightarrow Key-set Key-set

pre-halve(ks,sk,gk) \equiv card ks \geq 2
post-halve(ks,sk,gk) \equiv (sk \cup gk) = ks and (sk << gk) and
(card sks = card gks + 1)

where max: Nat-set \rightarrow Nat

pre-max(s) \equiv s \neq { }
post-max(s,r) \equiv r \in s and {r} \gg s

where \gg : Nat-set Nat-set \rightarrow Bool

s1 \gg s2 \equiv ($\forall e1 \in s1$) $\forall e2 \in s2$ (e1 \geq e2)

DELETEP

States: BP-tree

Type: Key \rightarrow

pre-DELETEP(l,k) \equiv k \in collect-keysp(TREEP(l))
post-DELETEP(l,k,l') \equiv ORDERP(l') = ORDERP(l) and
TREEP(l') = deleterp(DRDERP(l),TREEP(l),k)

where deleterp: Nat Nodep Key \rightarrow Nodep

pre-deleterp(m,n,k) \equiv k \in collect-keysp(n)
deleterp(m,n,k) \equiv let rn = deleterp(m,n,k)
if rn \in Inodep and sizep(rn) = 1
then TREEL(rn \times 1)
else rn

```

where deletep: Nat Nodep Key → Nodep
pre-deletep(m,n,k) ≡ k ∈ collect-keysp(n)
deletep(m,n,k) ≡ cases of n
  n ∈ Inodep: let it,ik,ck,cn = selectp(n,k)
               let rn = deletep(m,cn,k)
               if sizep(rn) ≥ minimum-sizep(m,rn)
               then let sn = << >, <cn>>
                   replace(n,sn,<< >, <m>>,<it>)
               else
                 let nn = neighbourp(n,it)
                 let sn = <<ck>, <cn,nn>>
                 if sizep(m) + sizep(nn) ≥
                   2 * minimum-sizep(m,rn)
                 then
                   replace(n,sn,redistributep(rn,nn,ck),ik,it)
                 else
                   replace(n,sn,mergep(rn,nn,ck),ik,it)
  n ∈ Tnodep: n \ {k}
end

```

```

where minimum-sizep: Nat Nodep → Nat
minimum-sizep(m,n) ≡ cases of n
  n ∈ Inodep: m+1
  n ∈ Tnodep: m
end

```

```

where neighbourp: Inodep Nat → Nodep
pre-neighbourp(n,i) ≡ i < len TREEL(n)
neighbourp(n,i) ≡ if i = len TREEL(n)
                  then TREEL(n)(i-1)
                  else TREEL(n)(i+1)

```

```

where redistributep: Nodep Nodep Key → Inodep
pre-redistributep(n1,n2,k) ≡ (((n1 ∈ Inodep and n2 ∈ Inodep) and
  (len KEYL(n1) + len KEYL(n2) > 2m-1) and
  (len TREEL(n1)=len KEYL(n1)+1 and len TREEL(n2)=len KEYL(n2)+1))
or (n1 ∈ Tnodep and n2 ∈ Tnodep and len TREEL(n1)+len TREEL(n2)
  > 2)) and
  ((collect-keysp(n1) <<= {k} << collect-keysp(n2)) or
  (collect-keysp(n2) <<= {k} << collect-keysp(n1)))

```

```

redistributep(n1,n2,k) ≡ cases of n1,n2
  n1,n2 ∈ Inodep: splitp o TREEL(mergep(n1,n2,k))X1)
  n1,n2 ∈ Tnodep: splitp o TREEL(mergep(n1,n2,k))X1)
end

```

```

where mergep: Nodep Nodep Key → Inodep
pre-mergep(n1,n2,k) ≡ (n1 ∈ Inodep and n2 ∈ Inodep) or
  (n1 ∈ Tnodep and n2 ∈ Tnodep)
mergep(n1,n2,k) ≡ cases of n1,n2
  n1,n2 ∈ Inodep:
    if collect-keysp(n1) << collect-keysp(n2)
    then << >, <KEYL(n1) || <k> || KEYL(n2), TREEL(n1) || TREEL(n2)>>
    else << >, <KEYL(n2) || <k> || KEYL(n1), TREEL(n2) || TREEL(n1)>>

```

```

  n1,n2 ∈ Tnodep: << >, <n1 + n2>>
end

```


4.4.4 Proofs of the correctness of Representation 2

4.4.4.1 Data Type Proofs

(B) Totality of the retrieve function

The rule to be proved is:

$$(1) (\forall t2 \in \text{BP-tree})(\text{invp}(t2) \Rightarrow \{\exists t1 \in \text{B-tree} \quad (t1 = \text{retrp}(t2) \text{ and } \text{inv}(\text{retrp}(t2)))\})$$

On substituting definitions, this becomes

$$(2) (\forall t2 \in \text{BP-tree})(\text{invrp}(\text{ORDERP}(t2), \text{TREEP}(t2)) \Rightarrow \\ \{\exists t1 \in \text{B-tree}(t1 = \langle \text{ORDERP}(t2), \text{retrp}(\text{TREEP}(t2)) \rangle \\ \text{and } \text{invr}(\text{ORDERP}(t2), \text{retrp}(\text{TREEP}(t2)))\})$$

Proof:

When *retrp* is applied to an element of *Inodep*, it is an identity function, and if it is applied to an element of *Inodep*, it relies only on the operator *elems* being applied to a list to be a total function. Since *TREEP*(*n*), where *n* \in *Inodep*, is always a list, and since the conditions of *invr* are weaker than those of *invrp*, it follows that the retrieve function *retrp* is total. This can be proved formally by structural induction on *TREEP*(*t2*).

(C) Adequacy

The rule to be proved is

$$(1) (\forall t1 \in \text{B-tree})(\text{inv}(t1) \Rightarrow \{\exists t2 \in \text{BP-tree}(\text{invp}(t2) \text{ and } \\ t1 = \text{retrp}(t2))\})$$

On substituting definitions, this becomes

$$(2) (\forall t1 \in \text{B-tree})(\text{invr}(\text{ORDER}(t1), \text{TREE}(t1)) \Rightarrow \\ \{\exists t2 \in \text{BP-tree} \\ (\text{invrp}(\text{ORDERP}(t2), \text{TREEP}(t2)) \text{ and } t1 = \langle \text{ORDERP}(t2), \text{retrp}(\text{TREEP}(t2)) \rangle)\})$$

Proof:

This is easiest shown by informal argument.

If *t1* \in *B-tree* and *TREE*(*t1*) \in *Inode* and *invr*(*ORDER*(*t1*), *TREE*(*t1*)) is true, then *t1* can be represented by *t2* = $\langle \text{ORDER}(t1), \text{TREE}(t1) \rangle$ and this clearly satisfies

$$\text{invrp}(\text{ORDERP}(t2), \text{TREEP}(t2)) \text{ and } t1 = \langle \text{ORDERP}(t2), \text{retrp}(\text{TREEP}(t2)) \rangle$$

For any *t1* \in *B-tree* such that *TREE*(*t1*) \in *Inode* and which satisfies *invr*(*ORDER*(*t1*), *TREE*(*t1*)) a representation *t2* can be produced as follows: *ORDERP*(*t2*) = *ORDER*(*t1*). Each *n* \in *Inode* of *t1* can be represented by the identical node in *t2* and will satisfy

$$\text{ORDERP}(t2) \ll \text{sizep}(n) \ll 2^{\text{ORDERP}(t2)}$$

because this invariant is identical to $ORDER(t1) \leftarrow size(n) \leftarrow 2 * ORDER(t1)$.

For each $n1 \in Node$ of $t1$, the elements of $n1$ can be laid out to form a list of nodes in ascending order of the key set that can be collected from each element (cf. keysets-are-orderedp) - this list forms $TREE1(n2)$. The maximum key in each of the key sets of the elements of $n1$ can be determined and all but the largest laid out in ascending order to form $KEYL(n2)$.

The fact that $t1$ satisfies $invt(ORDER(t1), TREE(t1))$ implies that each $n1 \in Node$ will satisfy $common-invt(n1)$ and $size-invt(ORDER(t1), n1)$ and this ensures that the corresponding node $n2 \in Nodep$ obtained as described will satisfy $common-invt(n2)$ and $size-invt(ORDERP(t2), n2)$ (cf. keysets-are-orderedp). Thus the resulting $t2 \in BP-TREE$ will satisfy $invt(ORDERP(t2), TREEP(t2))$ and $t1 = \langle ORDERP(t2), retrnp(TREEP(t2)) \rangle$.

4.4.4.2 Operation Proofs

Operation FINDP

(A) Preservation of the invariant

As FINDP is an identity operation on the BP-tree, it follows immediately that the invariant is preserved.

Proof that FINDP models FIND1

(D) Domains Rule

What must be shown is

$$(1) (\forall t \in BP-tree)(invt(t) \text{ and } pre-FINDP(retrnp(t), k) \Rightarrow pre-FINDP(t, k))$$

On substituting definitions this becomes

$$(2) (\forall t \in BP-tree)(invt(t) \text{ and } k \in collect-keys(retrnp(TREEP(t))) \Rightarrow k \in collect-keys(TREEP(t)))$$

Proof:

It follows that (2) is true from the fact that the structure of collect-keys is the same as that of collect-keysp (the extra key list present in each $n \in Nodep$ is ignored). Equation (2) can be proved formally to be true by structural induction on $TREEP(t)$.

(E) Results Rule

The rule to be proved is

- (1) $(\forall t \in \text{BP-tree}(\text{invp}(t)) \text{ and } \text{pre-FINDP}(t,k) \text{ and } \text{post-FINDP}(t,k',d) \Rightarrow \text{post-FIND}(\text{retrp}(t),k,\text{retrp}(t'),d))$

Expanding this gives

- (2) $(\forall t \in \text{BP-tree}(\text{invp}(t)) \text{ and } k \in \text{collect-keysp}(\text{TREE}(t)) \text{ and } t' = t \text{ and } d = \text{findp}(k,\text{TREEP}(t)) \Rightarrow \text{retrp}(t') = \text{retrp}(t) \text{ and } d = \text{find}(k,\text{retrnp}(\text{TREEP}(t)))$

It follows immediately from $t' = t$ that

$\text{retrp}(t') = \text{retrp}(t)$ so what must be shown is:

- (3) $(\forall t \in \text{BP-tree}(\text{invp}(t)) \text{ and } k \in \text{collect-keysp}(\text{TREEP}(t)) \Rightarrow \text{findp}(k,\text{TREEP}(t)) = \text{find}(k,\text{retrnp}(\text{TREEP}(t)))$

Proof:

Equation (3) can easily be shown to be true for elements of Inodep . When

$\text{TREEP}(t) \in \text{Inodep}$, it may be written:

$\text{let } n = \text{TREEP}(t)$

and writing only the relevant part of the invariant:

$\text{keys-are-orderedp}(n)$ and

$k \in \text{union}(\text{collect-keysp}(\text{TREEL}(n(i)) \mid 1 \leq i \leq \text{len } \text{TREEL}(n))) \Rightarrow$

$\text{findp}(k,\text{TREEL}(n)(\text{index}(k,\text{KEYL}(n)))) =$

$\text{find}(k,\text{select}(\{\text{retrnp}(n) \mid n \in \text{elems } \text{TREEL}(n)\},k))$

The invariant, keys-are-orderedp , guarantees that the function index causes that element of $\text{TREEL}(n)$ to be selected that corresponds to the "retrieved" element of $\text{TREEL}(n)$ that is chosen by the function select . So (4) can be seen to be true by the definitions of the functions index and select .

The results rule can be proved formally by structural induction on $\text{TREEP}(t)$ of a lemma:

$(\forall n \in \text{Inodep}(\text{invp}(n)) \Rightarrow$

$\text{let } \langle k,t \rangle = n \text{ in } t(\text{index}(k,k)) = \text{select}(\text{retrnp}(n)(k))$

Operation INSERTP

(A) Preservation of the invariant

The invariant of Representation 2 corresponds closely to that of Representation 1, the only real differences occurring in the definitions of `size-InvP` and `keysets-are-orderedP`, and being introduced by the fact that an `InodeP` contains a list of keys which is ordered, and the ordering of the `NodeP-list` is related to this key list, as is its length, whereas an `Inode` contains only a `Node-set`.

The structure of `INSERTP` is the same as that of `INSERT1`, and for all the functions used by `INSERT` there are corresponding functions of similar structure which are used by `INSERTP`. The only additional function which is used by `INSERTP` is the function `replace`, and this can be related to Representation 1 as follows:

Lemma 224.

$i, j, k \in \text{InodeP}$ and $n1 = \text{retrnp}(i)$ and $n2 = \text{retrnp}(j)$ and $n3 = \text{retrnp}(k)$
and $\text{is-subnode}(i, j)$ and $\text{within-bounds}(i, j, \text{collect-keys}(k))$ and
 $\text{keysets-are-orderedP}(i)$ and $\text{keysets-are-orderedP}(k)$ and
 $((\text{KEYL}(j) = < > \text{ and } nk = nt) \text{ or } (nk = \text{position}(\text{KEYL}(i), \text{KEYL}(j))))$
and $nt = \text{position}(\text{TREEL}(i), \text{TREEL}(j)) \Rightarrow$

$n2 \subseteq n1$ and $\text{retrnp}(\text{replace}(i, j, k, nk, nt)) = n1 - n2 \cup n3$ and
 $\text{keysets-are-ordered}(\text{retrnp}(\text{replace}(i, j, k, nk, nt)))$

where $\text{within-bounds} : \text{InodeP} \times \text{InodeP} \times \text{Key-set} \rightarrow \text{Bool}$
 $\text{pre-within-bounds}(i, j, ks) \equiv \text{is-subnode}(i, j)$
 $\text{within-bounds}(i, j, ks) \equiv \text{let } j = \text{position}(\text{TREEL}(i), \text{TREEL}(j))$
if $\text{KEYL}(j) = < >$
then let $k = j$
else
let
 $k = \text{position}(\text{TREEL}(i), \text{len } \text{TREEL}(j),$
 $\text{TREEL}(k))$
if $j = 1$ and $k = \text{len } \text{TREEL}(j)$

```

then true
else
if j=1
then ks << {KEYL(i)X0}
else
if k=len TREEL(i)
then ks >> {KEYL(i)Xj-1}
else
{KEYL(i)Xj-1} << ks <<= {KEYL(i)X0}

```

Therefore, a list of Lemmas for Insertion for Representation 2 can be drawn up, which matches Lemmas 1.1 to 1.11 of Representation 1 but refers to elements of lists rather than of sets. The only Lemmas which differ in form from those of Representation 1 will be those corresponding to L1.1 and L1.6; i.e. L2.1 and L2.6. These will become:

L21

$n \in \text{Inodep}$ and $\text{not-inp}(k,n)$ and $\text{size-invp}(m,n)$ and $n' = \text{insertnp}(m,n,k,d)$
 \implies
 $n' \in \text{Inodep}$ and $(\forall sn' \in \text{elems TREEL}(n')) \{sn' \in \text{Inodep}\}$ and
 $1 \leq \text{size}(n') \leq 2$ and $\text{len KEYL}(n') = \text{size}(n') - 1$

L26

$n \in \text{Inodep}$ and $\text{common-invp}(n)$ and $\text{not-inp}(k,n)$ and $\text{size-invp}(m,n)$ and
 $n' = \text{insertnp}(m,n,k,d) \implies$
 $n' \in \text{Inodep}$ and $1 \leq \text{size}(n') \leq 2$ and $\text{len KEYL}(n') = \text{size}(n') - 1$

With these two lemmas, the proof of Theorem 2.12 will follow the lines of the proof of Theorem 1.12 (to which it corresponds).

The only lemma for which the proof is somewhat different for Representation 2 is L2.7, and so it is proved below. L2.7 is concerned with the preservation of `keysets-are-orderedp` on insertion, which is satisfied provided the precondition of the following lemma is satisfied:

L2.25

$i1, i2, i3 \in \text{index}$ and $\text{is-subnode}(i1, i2)$ and
within $\text{bounds}(i1, i2, \text{collect-keys}(i3))$ and $\text{keysets-are-orderedp}(i1)$ and
 $\text{keysets-are-orderedp}(i3)$ and
 $(\forall i \in \{1 \dots \text{len TREEL}(i3)\}) (\text{depths}(\text{TREEL}(i3)(i)) = \text{depths}(i2))$ and
 $((\text{KEYL}(i2) < > \text{ and } nk = nt) \text{ or } (nk = \text{position}(\text{KEYL}(i2)(1), \text{KEYL}(i1))))$ and
 $nt = \text{position}(\text{TREEL}(i2)(1), \text{TREEL}(i1))) \implies$

$\text{let } r = \text{replace}(i1, i2, i3, nk, nt)$ in
 $\text{keysets-are-orderedp}(r)$ and
 $(\forall i \in \{1 \dots \text{len TREEL}(r)\}) (\text{depths}(\text{TREEL}(r)(i)) = \text{depths}(i1))$

The proof that `INSERTP` preserves the invariant `invp` therefore
models the proof that `INSERT!` preserves the invariant `inv`.

Proof of Lemma 2.7

To prove:

(1) $n \in \text{Nodep}$ and $\text{common-invp}(n)$ and $\text{not-inp}(k,n)$ and $n' = \text{insertnp}(m,n,k,d)$

\Rightarrow

$\text{keysets-are-orderedp}(n')$ and $\text{collect-keysp}(n') = \text{collect-keysp}(n) \cup \{k\}$

Proof: By structural induction on n

BASIS

(2) $n \in \text{Inodep}$

which follows immediately from the equivalent of L1.2 for Representation 2.

INDUCTIVE HYPOTHESIS $n \in \text{Inodep}$

assume

(3) $(\forall i \in \{1 \dots \text{len TREEL}(n)\})$

$(\text{let } sn = \text{TREEL}(nX_i) \text{ in}$

$\text{common-invp}(sn)$ and $\text{not-inp}(k,sn)$ and $sn' = \text{insertnp}(m,sn,k,d)$

\Rightarrow

$\text{keysets-are-orderedp}(sn')$ and $\text{collect-keysp}(sn') = \text{collect-keysp}(sn) \cup \{k\}$)

(4) $\text{common-invp}(n)$ and $\text{not-inp}(k,n)$

(Hypothesis)

(5) $(\forall i \in \{1 \dots \text{len TREEL}(n)\}) (\text{common-invp}(\text{TREEL}(nX_i))$

$\text{and not-inp}(k, \text{TREEL}(nX_i)))$

(by (4), common-invp and not-inp)

(6) $\text{let } i, k, ck, cn = \text{selectp}(n, k)$

(7) $cn \in \text{elems TREEL}(n)$

(selectp)

(8) $\text{common-invp}(cn)$ and $\text{not-inp}(k, cn)$

((5) and (7))

(9) $\text{let } cni = \text{insertnp}(m, cn, k, d)$

(10) $(\forall i \in \{1 \dots \text{len TREEL}(cni)\}) (\text{depths}(\text{TREEL}(cniX_i)) = \text{depths}(cn))$

(by L28 - equivalent of L1.8)

(11) $\text{keysets-are-orderedp}(cni)$

(by (7), (8) and (3))

(12) $\text{collect-keysp}(cni) = \text{collect-keysp}(cn) \cup \{k\}$

(by (7), (8) and (3))

(13) $\text{let } sn = \langle \langle \rangle, \langle cn \rangle \rangle$

(14) $\text{is-subnode}(n, sn)$

(by (7), (13) and is-subnode)

(15) $\text{collect-keyso}(sn) = \text{collect-keysp}(cn)$

(by (13) and collect-keysp)

(16) $\text{within-bounds}(n, sn, \{k\})$

(by (14) and selectp)

(17) $\text{within-bounds}(n, sn, \text{collect-keysp}(cni))$

(by (16), (12), and (15))

(18) $\text{let } rn = \text{replace}(n, sn, cni, i, i)$

(19) $\text{keysets-are-orderedp}(rn)$

((4), (10), (11), (14), (17), selectp , replace , L225)

(20) $\text{collect-keyso}(rn) = \text{collect-keysp}(n) \cup \{k\}$

((12), (15), replace L225)

(21) $n' = \text{if size}(rn) \leq 2^*m \text{ then } (rn) \text{ else splitp}(rn)$

The lemma holds in both cases, from the definitions of splitp , collect-keysp , $\text{keysets-are-orderedp}$, (19) and (20).

Proof that INSERTP models INSERT

(D) Domains Rule

What must be shown is

- (1) $(\forall t \in \text{BP-tree})(\text{invp}(t) \text{ and pre-INSERT}(t, k, d) \Rightarrow \text{pre-INSERTP}(t, k, d))$

On substituting definitions this becomes

- (2) $(\forall t \in \text{BP-tree})(\text{invp}(t) \text{ and } k \in \text{collect-keys}(\text{retrp}(t)) \Rightarrow k \in \text{collect-keys}(\text{TREEP}(t)))$

Proof:

This follows the proof of the domains rule for operation FINDP except that all occurrences of $k \in \dots$ are replaced by $k \in \dots$.

(E) Results Rule

The rule to be proved is

- (1) $(\forall t \in \text{BP-tree})(\text{invp}(t) \text{ and pre-INSERTP}(t, k, d) \text{ and post-INSERTP}(t, k, d, t') \Rightarrow \text{post-INSERT}(t, k, d, \text{retrp}(t')))$

Expanding this gives

- (2) $(\forall t \in \text{BP-tree})(\text{invp}(t) \text{ and } k \in \text{collect-keys}(\text{TREEP}(t)) \text{ and } \text{ORDERP}(t') = \text{ORDERP}(t) \text{ and } \text{TREEP}(t') = \text{insertrp}(\text{ORDERP}(t), \text{TREEP}(t), k, d) \Rightarrow \text{ORDERP}(t') = \text{ORDERP}(t) \text{ and } \text{TREEP}(t') = \text{insertrp}(\text{ORDERP}(t), \text{retrp}(\text{TREEP}(t)), k, d))$

Proof:

INSERTP is identical to INSERT1 when applied to Inodep (as a Inodep is identical to a Inode), and L2.24 and L2.25 show that INSERTP, when applied to an inodep, achieves the same effect as INSERT1 when applied to an Inode.

As INSERTP preserves the invariant invp on BP-trees, and it has been shown that FINDP models FIND1, this implies that the results rule (2) is true. It can be proved formally to be true by structural induction on TREEP(t).

Operation DELETEP

(A) Preservation of the invariant

As the invariant of Representation 2 is of a form similar to the invariant of Representation 1 and the structure of DELETEP corresponds to that of DELETE1, the proof that DELETEP preserves the invariant, $invp$, can be modelled on the proof that DELETE1 preserves the invariant, inv .

As in the case of insertion, a list of Lemmas for Deletion for Representation 2 can be drawn up, which matches L1.13 to Theorem 1.23 of Representation 1, but refers to elements of lists rather than sets. Also, the only additional function used by DELETEP is the function `replace`.

Because of the way in which `replace` can be related to Representation 1, (L2.24 and L2.25), it follows that the proof of Theorem 2.23 has the same structure as the proof of Theorem 1.23, to which it corresponds, and it can be modelled on the preceding proof.

Proof that DELETEP models DELETE1

(D) Domains Rule

What must be shown is

- (1) $(\forall t \in \text{BP-tree})(\text{invp}(t) \text{ and pre-DELETE1}(\text{retrp}(t),k) \Rightarrow \text{pre-DELETEP}(t,k))$

On substituting definitions this becomes

- (2) $(\forall t \in \text{BP-tree})(\text{invp}(t) \text{ and } k \in \text{collect-keys}(\text{retrnp}(\text{TREEP}(t))) \Rightarrow k \in \text{collect-keys}(\text{TREEP}(t)))$

Proof:

This proof is identical to the proof of the domains rule for operation FINDP.

(E) Results Rule

The rule to be proved is

- (1) $(\forall t \in \text{BP-tree})(\text{invp}(t) \text{ and pre-DELETEP}(t,k) \text{ and post-DELETEP}(t,k,t') \Rightarrow \text{post-DELETE1}(\text{retrp}(t),k,\text{retrp}(t')))$

Expanding this gives

- (2) $(\forall t \in \text{BP-tree})(\text{invp}(t) \text{ and } k \in \text{collect-keys}(\text{TREEP}(t)) \text{ and } \text{ORDERP}(t') = \text{ORDERP}(t) \text{ and } \text{TREEP}(t') = \text{deleter}(\text{ORDERP}(t),\text{TREEP}(t),k) \Rightarrow \text{ORDER}(t') = \text{ORDER}(t) \text{ and } \text{TREE}(t') = \text{deleter}(\text{ORDERP}(t),\text{retrnp}(\text{TREEP}(t)),k)$

Proof:

As a *Inodep* is identical to a *Inode*, *ODELETEP* is identical to *ODELETE1* when applied to a *Inodep*. It also follows, from L2.24 and L2.25, that *ODELETEP*, when applied to an *Inodep* has a similar effect to *ODELETE1* applied to an *Inode*.

ODELETEP has been shown to preserve the invariant, *invp*, on *BP-trees* and as it has been shown that *FINDP* models *FIND1*, it follows that (2) is true. It can also be proved to be true formally by using structural induction.

4.5 Realization

This stage consists of the program code corresponding to the FIND, INSERT and DELETE operations for B⁺-trees. The system has been implemented in PASCAL, using the UCSD system on a DEC LSI-11.

The code for FIND and INSERT contains weakest pre-condition type assertions as an indication of its correctness. Assertions could be added to the code for DELETE in a similar manner.

The mapping of the B⁺-tree onto linear storage is embodied in the pointer type variables of PASCAL and a form of has-no-loops (cf. section 3.3) for B⁺-trees corresponding to that for binary trees should be true of these pointers.

When the height of the tree increases, a new node is allocated by using the new Command of PASCAL. As the version of the UCSD system which was used (Version II.0) actually implements the heap as a stack, a corresponding release command could not be used to free nodes if the height of the tree should decrease, because the system cannot be relied upon to discard nodes in the order in which they were acquired.

As a function in PASCAL can only return a result which is of a simple type, the function INSERTN could not return an actual node. A globally allocated node, r_node, is therefore used to hold intermediate results produced by INSERTN, and the function returns a pointer to this node at each instance.

The layout of the node was chosen following the form of representation used in [7], as this shortens sections of the program code considerably. (The implementation was originally

based on a node containing separately defined lists of different lengths so an actual comparison of code lengths was made). The additional pointer in the case of an `lnode` is stored in the node variable , `pmml`. In the case of a `Tnode`, this variable will be used to hold the sequential link to the following `Tnode` when the `NEXT` operation is implemented.

The program listing appears in the Appendix.

5. CONCLUSION

The most important aspect in the development of the representation of abstract mappings as B^+ -trees, was the choice of correct development stages. It was crucial that each stage should capture just the right aspects and properties to allow the steps between successive stages to be made easily and smoothly. The correct choice of representation at each level makes the refinement process a natural progression.

The main difficulty in the refinement proofs was in proving that the operations for insertion and deletion preserve the invariant. This problem was solved by decomposing the rule to be proved into a series of lemmas in each case, from which the preservation of the invariant follows logically. The separation of lemmas into a clear sequence was not an easy task, but the resulting sets of lemmas greatly simplify the proofs of validity.

Another problem is how to relate error-handling in the actual program code to pre-conditions which occur in the specification. This could be done by appropriate sets of assertions in the program code. However in this development a comment has simply been inserted in the code where such checking occurs.

In conclusion, possible future extensions to this work should

be considered. Briefly, a few of these are:

The implementation of the NEXT operation, for which provision has already been made in the node record layout.

The implementation of B⁺-trees using disk files and disk file addresses rather than storage and pointer variables and the extension of this to provide a separate file access package for general use (the difficulties of the strong typing of PASCAL have to be overcome to achieve this extension).

After the above two extensions the implementation of recovery procedures becomes possible.

Finally, allowing features such as key compression and multi-user access are also possible future developments.