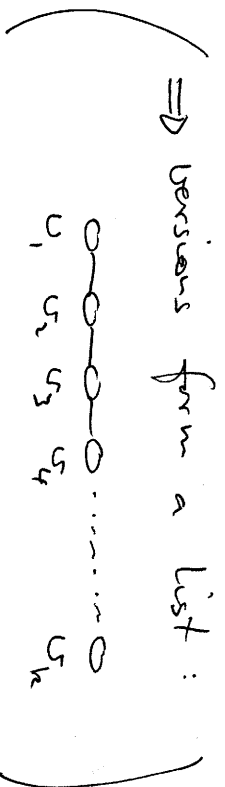


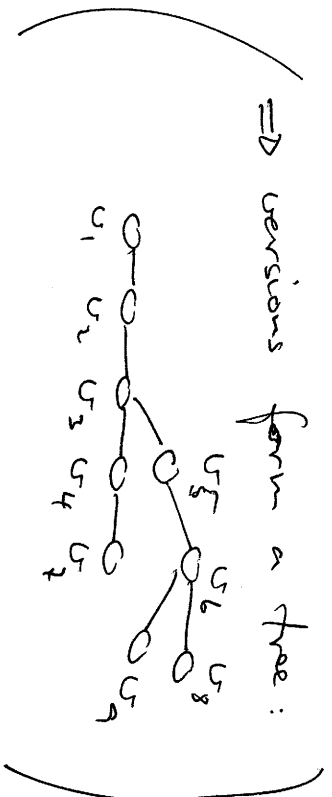
Persistent B-trees

Persistence : Each update gives a new version of tree
Give access to all versions (created so far).
(so similar in spirit to version control).

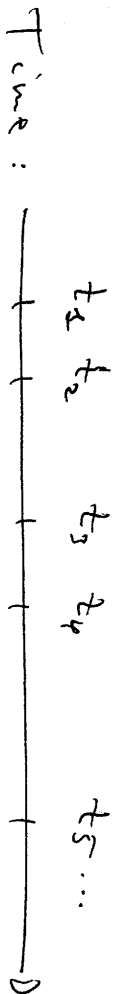
Partial persistence : Search in all versions
update in newest only



Full persistence : Search in all versions
Update \rightarrow _____



- + Naive solution : copy structure (entire) at each update
 - + We want more efficient solutions (in terms of time and space).
- Generic solutions for pointer-based data structures (in RPN model) developed by Pincoll et al (1989)
- TODAY : I/O-efficient partial persistent B-trees



Version: 1 2 3 4 5

Version i exists in time interval $[t_i; t_{i+1}]$
 (I.e., i th update starts at time t_i)

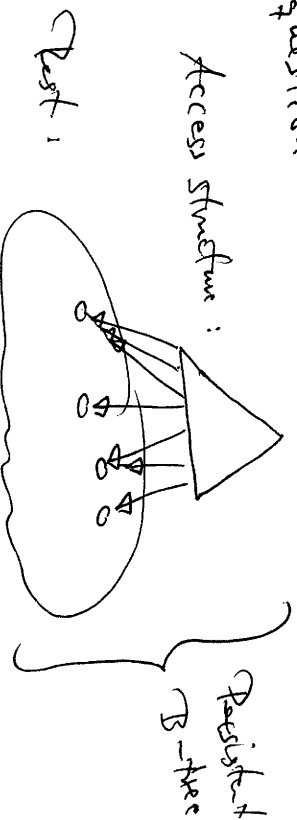
Using (t_i, i) as (key, value) - pairs stored in an ordered dictionary, we can find the ~~correct~~ number of the version existing at any time t (by predecessor search on t).

ie, non-persistent

Below, we will use a (standard) B-tree as this ordered dictionary, and will call it the access structure.

Except the values will not be version number, but a pointer into the rest of our (persistent B-tree) structure - namely to the correct place to start the search in the version in question

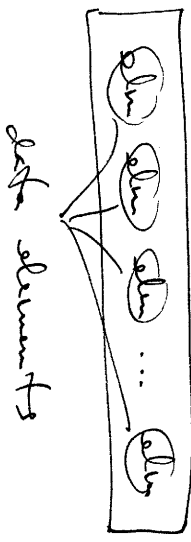
Access structure:



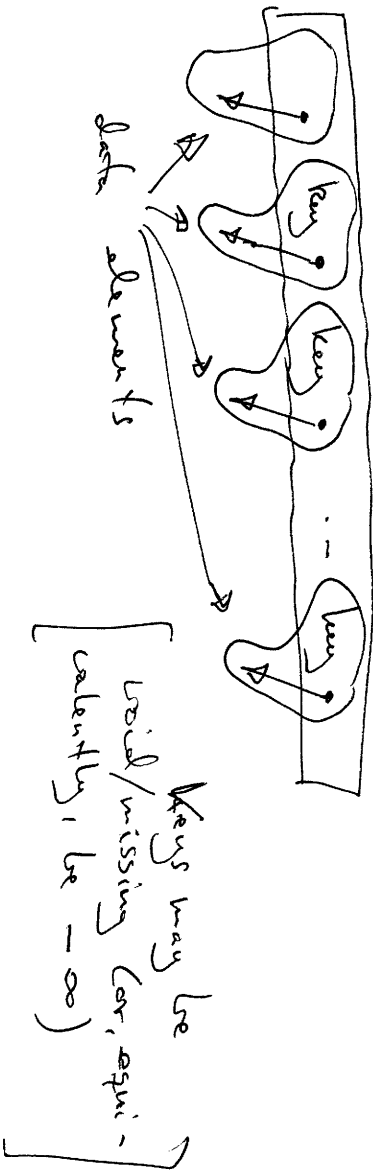
The rest : Let $l \geq 16$ be a parameter.

The rest of our structure consists of nodes containing data elements up to l .

For leaf nodes, the data elements are the elements stored ~~in~~ in our persistent B-tree (i.e., are (key, value) - pairs)



For internal nodes, the data elements are pointers (i.e., (key, node-pointer) - pairs).



Hence, "the rest" forms a graph.

Unlike a (standard) B-tree, the graph is not restricted to be a (unithin) tree. In fact, it will be a DAG.

We choose l such that a disk block can hold l data elements (l is B).

Also different from \mathbb{B} -trees: each data element has a time stamp.

A time stamp is an interval of the form

$$i) [t_i; t_j[$$

or

$$ii)]t_i; \infty[$$

for t_i, t_j update times. It states the versions for which the data element is valid.

Let t be the current time (ie, $t \in [t_i; \infty[$ for i the (currently) last version).

A data element is said to be alive if t is contained in its time stamp, and dead otherwise. (In fact, it will turn out that alive data elements have the $]t_i; \infty[$ form of the time stamp, ~~and~~ and dead elements the other form). from a note

No data element will ever get deleted. Time-stamps may change, however (by setting $\infty \rightarrow t_r$) (of alive nodes). We call \uparrow this closing the time-stamp.

A data element is valid at (some) time t if t is contained in its timestamp (so alive = valid at current time).

A node is valid if it contains any data elements valid at time t .

The main invariant maintained by the "rest" part of our structure is:

INVARIANT 1: For any t , the nodes ~~are~~ ^{valid} at t and the data elements valid at t together constitute a $(\frac{1}{4}, \frac{1}{4})$ -tree with leaf parameter $\frac{1}{4}$, and this tree contains in the leaf nodes exactly the elements stored in the version existing at time t .

In our access structure, the pointer stored as the value attached to t_i points to the root of the tree of the invariant for $t = t_i$.

Searches : The operation

$\text{Search}(N_t, t)$

that searches for the element of leaf l in version existing at time t , is now straight-forward :

Search in access structure with t (predecessor search).

Search as usual in (a, b) -trees with leaf l , starting at root pointed to by pointer returned by first search.

In details:

Each time a node is read, extract the data plus. valid at time t , sort these, view the result as an (a, b) -tree node and proceed

Each node access is one I/O. Hence, entire search takes $O(\log(I)) + O(\log(N_t))$

$$\cancel{\left(\sum_{t=1}^I O(\log(I + N_t)) \right)} = O(\log(I))$$

where I is the current number of versions (= number of updates so far) and N_t is the number of obs. in version existing at time t .

$I \geq N_t$ if we start with empty tree (assumed here)

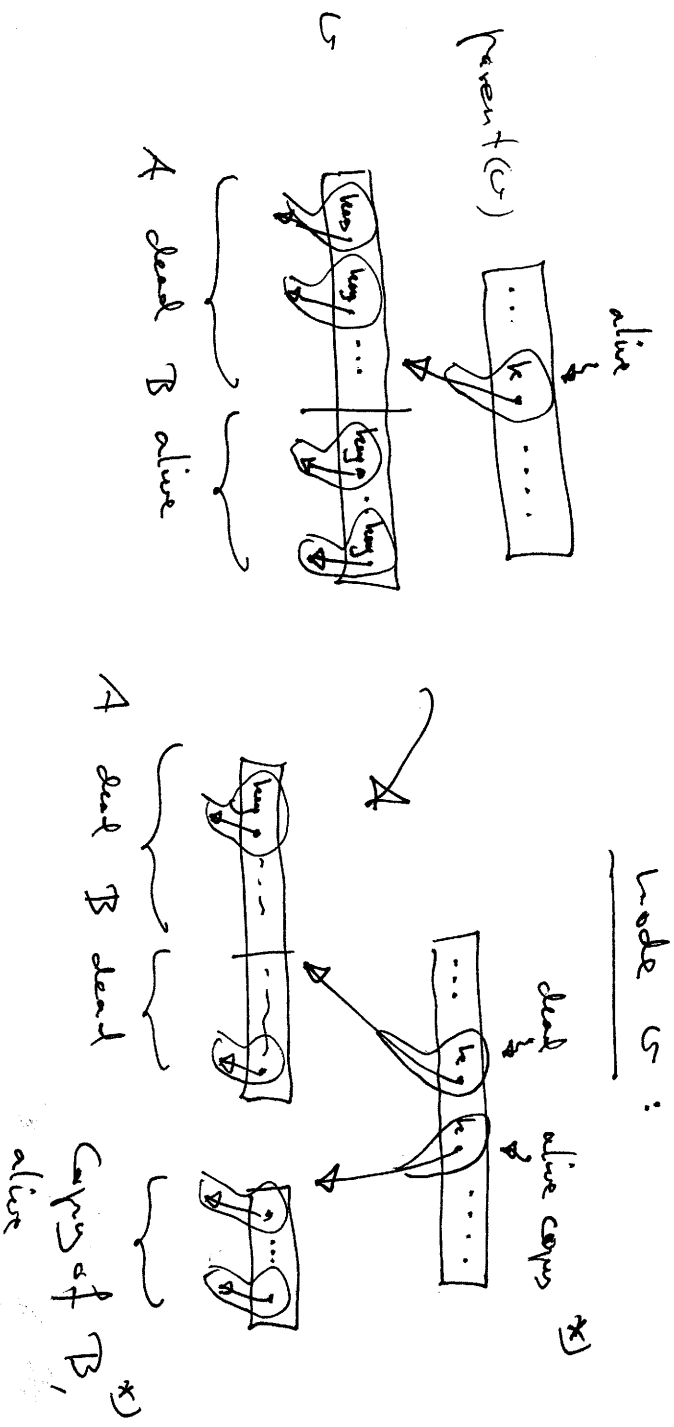
Notes : i) The above is I/O-efficient, but not CPU-optimal. (For this, nodes should contain persistent binary search trees, made using methods from Driscoll et al., which is out of scope here), see eg. to the sorting.

ii) If searches are based on ~~versions~~ version numbers (instead of time t) directly, the access structure can be simply an array of pointers (using O(1) I/O's to look up).

Updates

[Note: for updates $t = \text{current time}$ (as partial persistence)]

A basic operation will be the Version Split of a



+ Current time is t_i (if the number of the new version made by the update)

+ Existing dead timestamps are left untouched

+ The existing alive timestamps are made dead by changing $[t_i; \infty]$ to $[t_i, t_i]$

+ New alive copies get timestamp $[t_i, \infty]$ \uparrow (*)

+ Operation will only be included on search path for update position in current version. Hence the data element of the parent is alive (is assumed in drawing above).

Observation: A version split does not violate the invariant INV 1.

(Proof: only two previous / three new nodes need to be considered. Observe that for any t , what is valid at time t is not changed).

To obtain a good space analysis in the end, we introduce the following invariant:

INV 2: New nodes ~~are~~ created with $\frac{3}{8}L \leq \# \text{ data elements} \leq \frac{7}{8}L$

Also recall the already introduced constraint:

INV 3: No node ever contains more than L data elements.

We now proceed to define the update operations.

Insert :

Increment version counter i
 Log the update time t_i
 Search in persistent B-tree for insert-
 tion point (using time t_{i-1}).
 Insert data element in found leaf node
 with time stamp $[t_i; \infty[$

If needed : rebalance.

↙ ↘
 (leaf's below)

Needed :

INV 3 may be violated ($l+1$ data elements
 in leaf)
 INV 2 may also be violated for $t_i \leq t$
 (a B-tree overflow) - but then the above
 violation also occurs. So enough to check
 that.

Rebalance :

~~Make~~ Version Split ~~on~~ leaf.

IF INV 2 is violated by overflow ($> \frac{7}{8}l$
 data elem. in new node (call alive)):
 make a standard split (B-tree split) on
 new node (inserting a new data elem. in
 parent (with timestamp $[t_i; \infty[$)).
 New nodes will have at least

$$\lfloor \frac{7}{8}l + 1 \rfloor \geq \frac{7}{16}l > \frac{3}{8}l$$

~~data elem.~~
data elem.

and at most

$$\lceil \frac{l+1}{2} \rceil \leq \frac{l+1}{2} = l(\frac{1}{2} + \frac{1}{2l})$$

$$\leq l(\frac{1}{2} + \frac{1}{16}) < \frac{7}{8}l$$

data elem.

If FNV2 is violated by underflow

($< \frac{3}{8}L$ data elems. in new node (all alive)):

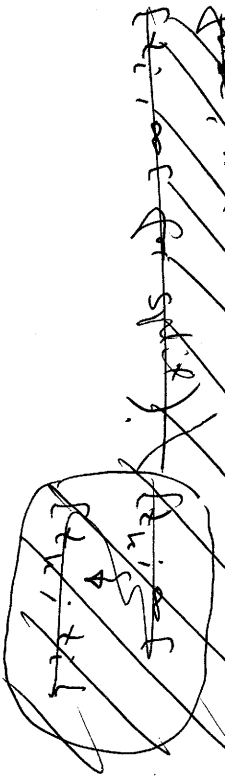
Find a sibling in t_i -B-tree (via alive data elems in parent).

Make a Version Split on this. (this gives between t_L and L alive data elems. due to sibling being a valid B-tree node at time t_i (FNV 2)).

Make a standard B-tree fuse between the two new nodes.

If result has $> \frac{3}{8}L$ data elems, make a standard B-tree split.

Update parent. ~~Change a timestamp for the fuse, add new data elem with timestamp~~



After the fuse, the new node has at least

$$\frac{1}{4}L + \frac{1}{4}L = \frac{1}{2}L > \frac{3}{8}L$$

data elems. With no split, it has at most $\frac{3}{8}L$.

After a split, the two new nodes have at least $\lfloor \frac{7}{8}L + 1 \rfloor > \frac{3}{8}L$ data elems (see previous pag). ~~and at~~

$$\text{most } \left\lceil \frac{\frac{3}{8}L - 1 + L}{2} \right\rceil \leq L \left(\frac{\frac{3}{8} + \frac{1}{2}}{2} \right) = \frac{11}{16}L < \frac{7}{8}L.$$

The two version splits introduced two new data elements (with timestamps $t_i; \infty$).

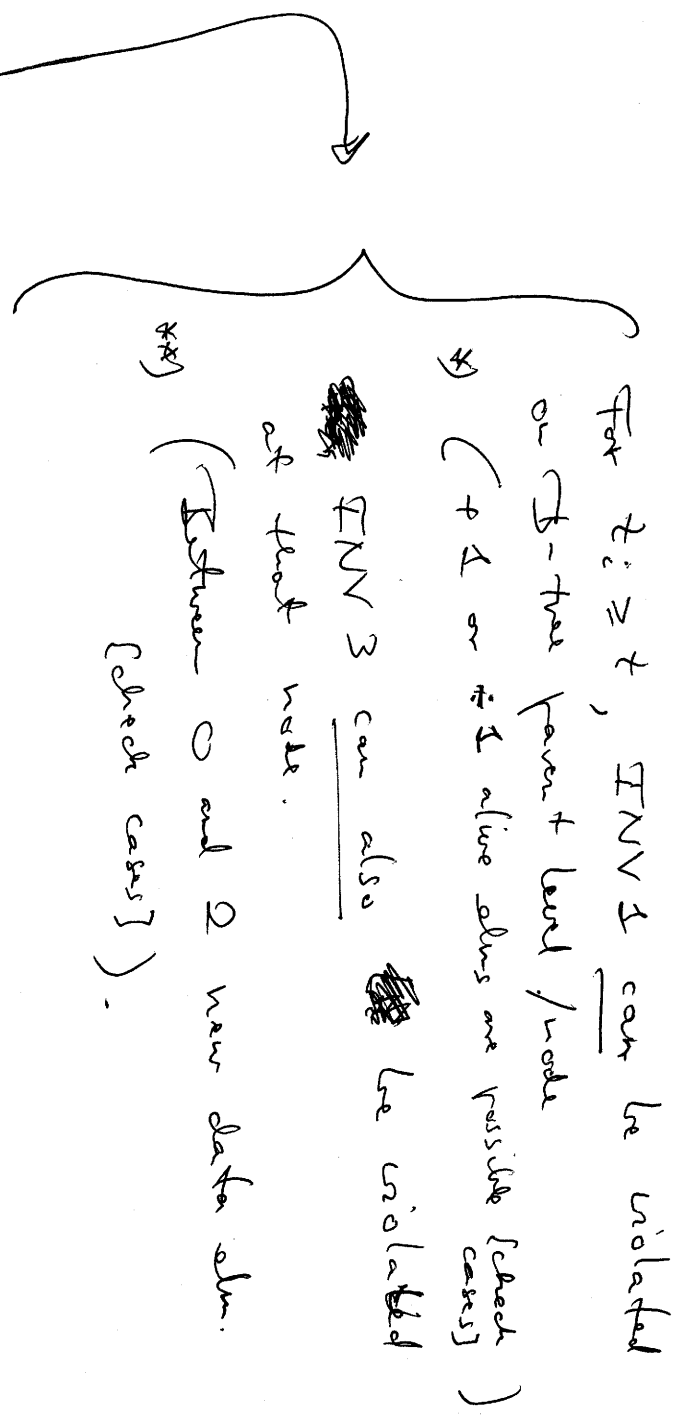
~~Now adjust these:~~
~~two data elements~~

Fuse: remove one
 Split: add one again.

(Fused nodes are valid B-tree nodes at time t_i)

Obscure: ^{Now,} No violations of INV 2 on current level (B-tree level) [i.e. on newly produced nodes]. See calculations above.

No violations of INV 1 (only $t \geq t_i$ need to be considered) or INV 3 (as levels of INV 2 are stronger) on the current B-tree level.



For $t_i \geq t$, INV 1 can be violated on B-tree parent level / node

*) ($t \leq$ or $t_i \leq$ alive plus one possible [check cases])

~~INV 3~~ can also ~~be~~ be violated at least node.

(Estimate 0 and 2 new data elem. [check cases]).

Its fixed by recursion on parent, starting with a Version Split on this.

Except that some node size counts may be large / smaller by one (which doesn't invalidate data and conclusions of calculations [check]), the same analysis applies.

Recursion may propagate to rest of t_i - B-tree, creating a new t_i - rest.

Therefore (at end of insertion) update access structure [insert (t_i , pointer to rest)]
 ↙
 new or not.

Deletion is very similar. Start by closing a file story in leaf. Analysis and actions are very similar to insert (astrakally, fewer cases)

Bottom: note the entire analysis based on a change at current level [initially leaf level] of $\pm I$ line data element change $0 \leq x \leq 2$ data element change.

(This analysis of recursive cases and of deletion is built in from beginning.)

Note:

During the recursive rebalancing, `rebal` at a given level is triggered by a) isolation of `INV 1` for t_i (current, alive tree), where 3-tree nodes may be isolated, or ii) isolation of `INV 3` (too many data elem. in node).

Denote these changes steps.

- 1) New nodes contain between $\frac{3}{8}L$ and $\frac{7}{8}L$ alive data elements (and nothing else). So needs $\frac{1}{8}L$ new data elements or $\frac{1}{8}L$ changes of number of alive data elements before they trigger updates.

A node can only trigger once (then made dead by a version split, and only alive nodes (ie, nodes in current tree) ~~are~~ are part of updates.

- 1) Each update gives one step on level 0 (in current tree)
- 2) Each `rebal` at level i gives at most two steps on level $i+1$ [See (b) and (c)] some pages above

— the relevant count is max (not sum) of the two types of steps.]

(numeric value of)

With I updates (= # versions) :

$\Rightarrow I$ steps on level 0.

$\Rightarrow \leq I / \frac{1}{b} = I \cdot \frac{8}{b}$ rel. ~~on~~ level 0

$\Rightarrow \leq I \cdot \frac{8}{b} \cdot 2$ steps on level 1

$\Rightarrow \leq I \cdot \frac{8}{b} \cdot 2 / \frac{1}{b} = I \cdot \left(\frac{8}{b}\right)^2 \cdot 2$ rel. ~~on~~ level 1

$\Rightarrow \leq I \cdot \left(\frac{8 \cdot 2}{b}\right)^2$ steps on level 2

$\Rightarrow \leq I \cdot \left(\frac{8 \cdot 2^2}{b}\right)^2 / \frac{1}{b} = I \cdot \left(\frac{8}{b}\right)^3 \cdot 2^2$ rel. on level 2

\vdots

Total : $\leq I \cdot \frac{8}{b} \cdot \sum_{i=0}^{\infty} \left(\frac{8 \cdot 2^i}{b}\right)^i$

$$= I \cdot \frac{8}{b} \cdot O(1) \quad \text{as } b > 16$$

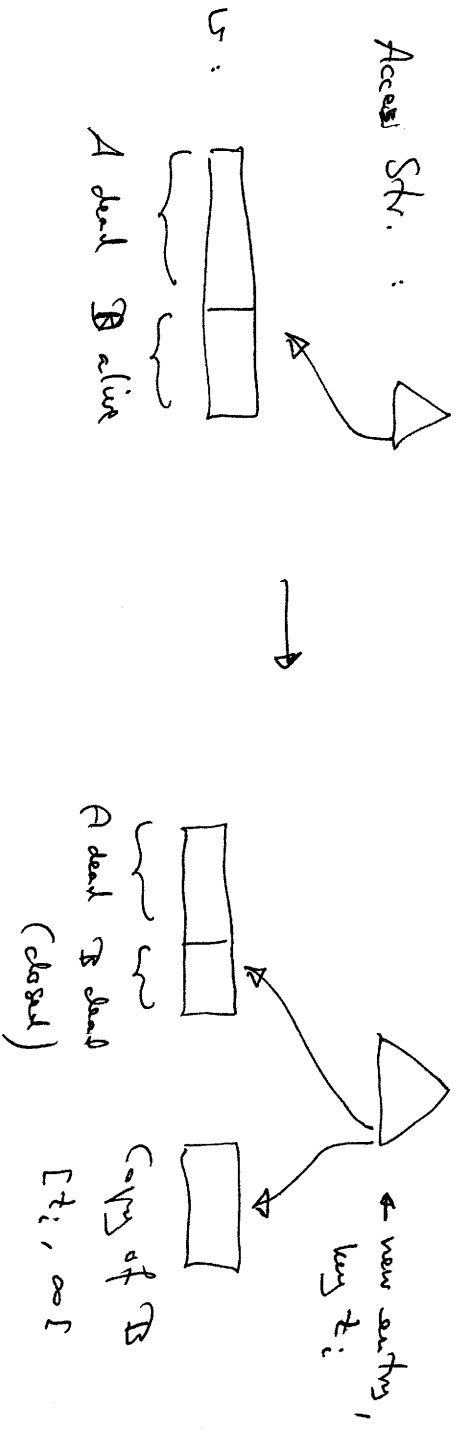
Time: Amortized $O\left(\frac{1}{b}\right)$ rel. (= I/b) for update.

Space: Each rel. \Rightarrow at most 2 new nodes. So total space $\leq 2 \cdot$ total rel.

$$= O\left(\frac{I}{b}\right)$$

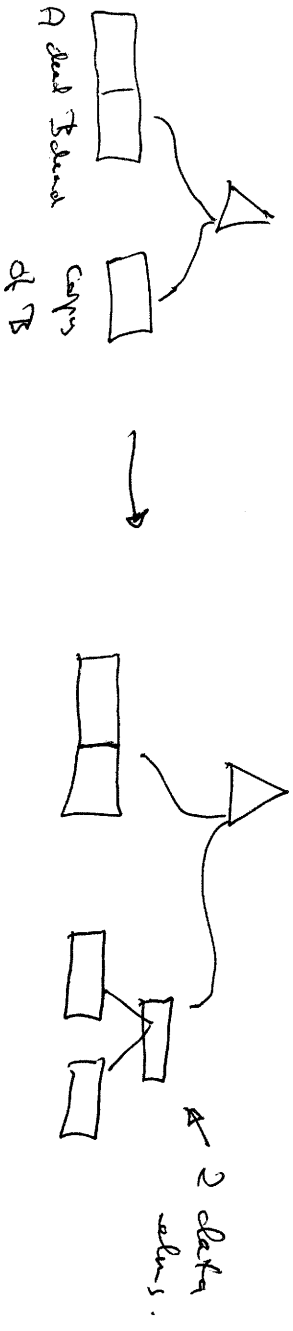
Special conditions at root :

Version Split at root :



Note that under flow ($< \frac{7}{8}L$) cannot be handled (there is no sibling to fuse with).

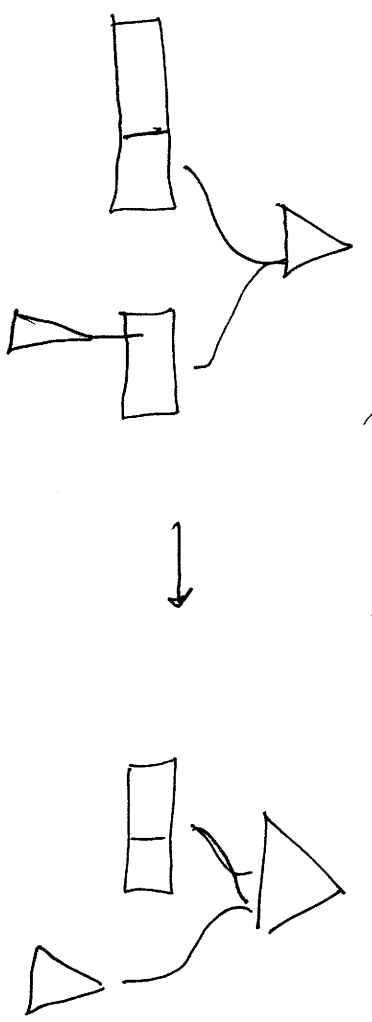
Also, the over flow case ($> \frac{7}{8}L$) has a split which at root gives a node with two data slots :



So INV 2 has to have the following addition :

INV2' : If new node is root of current tree, then $2 \leq \# \text{ data slots} \leq \frac{7}{8}L$

However,
 If the underflow of new root (after
 Version Split) is 1, we do as follows :
 (data always present)



I.e. we remove the node.

(Note that underflow cannot be 0 data always present, as the current tree was a valid B-tree before update, and at least 1 child
 (So root here has ≥ 2 children)

was removed by the recursive rebalancing marking the root.)

The changed EM 2 affects the analysis of time and space:

The analysis above ~~is~~ really ~~not~~ ^{does not} cover

rebalancing of current not triggered by water flow.

Such triggering could happen after 1 (not 1/2 b) steps in the node [So a) does not hold?]

However, such relab. also does not imply steps on the level above.

So analysis on page (15) is valid for the rest of the relab. ops.

(root relab. triggered by water flow)

~~As~~ such triggering happens at levels ≥ 1

(Really, a root at level 0 is B-free must be allowed down to 0 (not 2) steps), and requires $\Omega(1)$ steps, line 4 on page (15) shows that such relab. is $O(I/b)$.

Hence, time and space bounds are still valid. of page (15), bottom.