# PMA LAYOUT

## DISK

^ 2x pages



2MB

DATA    4MB    ~ DATA    . . .

## MEMORY *

2MB    4MB



~ DATA ~ ~ ~ ~ ~ . . .

↑ arena base ** — $2^{45}$ = 0x2000.0000.0000 — above exe/heap

* VIRTUAL MEMORY

** top — $2^{45}$ + $2^{46}$ = 0x6000.0000.0000 — below shared libs/stack

# SINGLE TREE MAX SIZE

$$\frac{2^{46} \text{ (arena bytes)}}{2^{14} \text{ (page bytes)}} = 2^{32} \text{ pages}$$

$$\frac{2^{32} \text{ pages}}{2^{10} \text{ entries/leaf}} = 2^{22} \text{ leaves}$$

$$\sim 2^{23} \text{ nodes}$$

$$2^{23} \text{ nodes} \cdot 2^{14} \text{ bytes/node} = 2^{37} \text{ bytes} = 2^{7} GB = 128 GB$$

# NODE LAYOUTS

## ROOT NODE

| | |
|---|---|
| 4 bytes - version | 4 |
| 4 bytes - checksum | 8 |
| 8 x 4 bytes - block bases  ⎱ disk offsets of | 16 |
| | ⎰ spaces for nodes | |
| 1 byte - block base extent (0-8) | 17 |
| 1 byte - tree depth | 18    = 69 words |
| 1 byte - dirty flag | 19    out of 4K |
| 1 byte - padding | 20 ⎰ 2 bytes entry count |
| 256 bytes - dirty bitmap | 276 ⎱ 2 bytes padding |
| M F M F M ...F M   remainder | |

↑          ↑

Ø       MAX-OFFSET

M and F  32-bit (4 byte)  words

# NODE LAYOUTS

## ~~ROOT~~ NODE

4 bytes - version                                                    4
4 bytes - checksum                                                   8
8 x 4 bytes - block bases ( disk offsets of                         16
1 byte - block base extent (0-8)                                    17
1 byte - tree depth                                                 18    = 64 words
1 byte - dirty flag                                                 19         out of 4K
1 byte - padding                                                    20
256 bytes - dirty bitmap                                        ~~256~~ 256
m F M F m ...F M   remainder

( disk offsets &
spaces for nodes

M and F  32-bit (4 byte)  words

128 GB — b[7]     (128 GB)

32 GB — b[6]     (160 GB)

8 GB — b[5]     (168 GB)

2 GB — b[4]     (170 GB)

512 MB — b[3]     (170 GB + 512 MB)

128 MB — b[2]     (170 GB + 640 MB)

32 MB — b[1]     (170 GB + 672 MB)

8 MB — b[0]     (170 GB + 680 MB)

2 MB — fixed     (170 GB + 682 MB)

$$\frac{\sim 171\ GB}{64\ TB - \sim 171\ GB} =$$

$$\frac{171 \cdot 2^{30}}{64 \cdot 2^{40} - 171 \cdot 2^{30}}$$

$$\approx \frac{171 \cdot 2^{30}}{63 \cdot 2^{40}}$$

$$\approx \frac{171}{63} \cdot \frac{1}{2^{10}}$$

$$\approx 4 \cdot \frac{1}{2^{10}}$$

$$\approx \frac{1}{2^8} \quad \frac{1}{256}$$

$< 1\%$ overhead

$$2^{22} \quad \text{leaves} \quad - \quad 64 \text{ GB} \quad \cdot 2 \quad = \quad 128 \text{ GB} \quad (128 \text{ GB slot})$$
$$+$$
$$2^{12} \quad \text{d-1} \quad \sim \quad 64 \text{ MB} \quad \cdot 2 \quad = \quad 128 \text{ MB} \quad (128 \text{ MB slot})$$
$$+$$
$$2^{2} \quad \text{d-2} \quad - \quad 64 \text{ KB} \quad \cdot 2 \quad = \quad 128 \text{ KB} \quad (8 \text{ MB slot})$$
$$+$$
$$1 \quad \text{root} \quad - \quad 16 \text{ KB} \quad \cdot 2 \quad = \quad 32 \text{ KB} \quad (2 \text{ MB slot})$$

Enough room to fully dirty the tree

(Why not sparse files?)

fallocate(2)
↙ used for sparse
  file hole punching

unclear interaction between
fallocate(), mmap(), msync()

# DELETION

# Deletion without Rebalancing in Multiway Search Trees*

Siddhartha Sen[1] and Robert E. Tarjan[1,2]

[1] Princeton University
{sssix,ret}@cs.princeton.edu
[2] HP Laboratories, Palo Alto CA 94304

Just delete the entry;
the leaf if now empty,
recursively the nodes.

**Abstract.** Many database systems that use a $B^+$ tree as the underlying data structure do not do rebalancing on deletion. This means that a bad sequence of deletions can create a very unbalanced tree. Yet such databases perform well in practice. Avoidance of rebalancing on deletion has been justified empirically and by average-case analysis, but to our knowledge no worst-case analysis has been done. We do such an analysis. We show that the tree height remains logarithmic in the number of insertions, independent of the number of deletions. Furthermore the amortized time for an insertion or deletion, excluding the search time, is $O(1)$, and nodes are modified by insertions and deletions with a frequency that is exponentially small in their height. The latter results do not hold for standard $B^+$ trees. By adding periodic rebuilding of the tree, we obtain a data structure that is theoretically superior to standard $B^+$ trees in many ways. We conclude that rebalancing on deletion can be considered harmful.

# FREE LISTS

$$0x6000.0000.0000$$
$$= 0x2000.0000.0000 = 2^{45}$$
$$+ 0x4000.0000.0000 = 2^{46}$$

- memory
- node disk
- pending node disk
- data disk
- pending data disk

} in data memory

```
            128 TB
       64       64 TB
0      ↓        ↓        ↓
↓      |    |———————|    | |
↑               ↑
exe + heap      stack +
                shared objs
```

<span style="color:green">Unify pending just prior to snapshot:
allow no new dirtying from tree list
until snapshot succeeds. Process
crashes if snapshot fails</span>

<span style="color:red">* what about online snapshotting?</span>

1   - 2 MB only
      - init / find / insert / delete
                 (split)

2 - memory management (take disk offsets)
     - data arena            - malloc / free
     - memory free lists     - free can require
                                     huge or allocation

3 - disk persistence
     - disk free lists
     - dirtying / cow
     - msync / cleaning        5 - range freeing
     - restoring /

4 - node partition striping

# insert into empty range

O O U32:MAX

| O | >5 | O |
|---|----|---|
| 0 | 2 3 | U32:MAX |

# insert into full range

| 89 |
15    21

16  11 2 1  18

| 89 | 11 2 1 9 2 1 |
15   16    18    21