# AVL 树 → 红黑树问题

聂浩　无 31　2013011280

2015 年 12 月 26 日

## 1　实验内容

在 Windows 的虚拟内存管理中，将 VAD 组织成 AVL 树。VAD 树是一种平衡二叉树。

红黑树也是一种自平衡二叉查找树，在 Linux 2.6 及其以后版本的内核中，采用红黑树来维护内存块。

请尝试参考 Linux 源代码将 WRK 源代码中的 VAD 树由 AVL 树替换成红黑树。

## 2　实验思路

红黑树与 AVL 的查找基本是一致的，不需要太多的修改。

插入与删除函数通过对 linux(v4.4-rc6) 与 WRK1.2 的分析。利用用 linux 的红黑树代码替换 WRK 中的 AVL 管理。

### a)　Linux 中的红黑树

Linux 中，虚拟内存管理的 VAD 由红黑树实现。

红黑树是一种在插入或删除节点是都需要维持平衡的二叉查找树，且每个节点都具有颜色属性:

1. 一个结点要么是红色的，要么是黑色的。

2. 根结点是黑色的。

3. 如果一个结点是红色的，那么它的子结点必须是黑色的，也就是说在沿着从根结点出发的任何路径上都不会出现两个连续的红色结点。

4. 从一个结点到一个 NULL 指针的每条路径上必须包含相同数目的黑色结点。

本次实验使用 Linux Kernel 4.4-rc6 版本。

linux 内核源代码中，红黑树的定义在三个地方，分别是

- include/linux/rbtree.h

- include/linux/rbtree_augmented.h

- lib/rbtree.c

其中红黑树节点的定义为：

```
struct rb_node {
    unsigned long  __rb_parent_color;
    struct rb_node *rb_right;
    struct rb_node *rb_left;
} __attribute__((aligned(sizeof(long))));
    /* The alignment might seem pointless, but allegedly CRIS needs
        it */

struct rb_root {
    struct rb_node *rb_node;
};
```

有意思的是，这里使用了 \_\_attribute\_\_((aligned(sizeof(long)))), 使得结构体进行 4 字节大小对齐（32 位系统），因此地址最低两位始终是 00，linux 使用其中的最低一位表示红黑树节点的颜色。

基于此 linux 定义了一些基本操作（在 rbtree_augmented.c 中）：

```
\\在 rbtree.h 中
#define rb_parent(r)    ((struct rb_node *)((r)->__rb_parent_color & ~3))
\\在 rbtree_augmented.c
#define RB_RED        0
#define RB_BLACK      1

#define __rb_parent(pc)    ((struct rb_node *)(pc & ~3))

#define __rb_color(pc)      ((pc) & 1)
#define __rb_is_black(pc)   __rb_color(pc)
#define __rb_is_red(pc)     (!__rb_color(pc))
#define rb_color(rb)        __rb_color((rb)->__rb_parent_color)
#define rb_is_red(rb)       __rb_is_red((rb)->__rb_parent_color)
#define rb_is_black(rb)     __rb_is_black((rb)->__rb_parent_color)

static inline void rb_set_parent(struct rb_node *rb, struct rb_node *
    p)
{
    rb->__rb_parent_color = rb_color(rb) | (unsigned long)p;
}
// 设置父节点的同时设置自己的颜色
```

```
21  static inline void rb_set_parent_color(struct rb_node *rb,
22                          struct rb_node *p, int color)
23  {
24      rb->__rb_parent_color = (unsigned long)p | color;
25  }
26
27  static inline void
28  __rb_change_child(struct rb_node *old, struct rb_node *new,
29          struct rb_node *parent, struct rb_root *root)
30  {
31      if (parent) {
32          if (parent->rb_left == old)
33              WRITE_ONCE(parent->rb_left, new);
34              \\WRITE_ONCE(a,b)=>a=b;
35          else
36              WRITE_ONCE(parent->rb_right, new);
37      } else
38          WRITE_ONCE(root->rb_node, new);
39  }
40  \\在rbtree.c中，为旋转做准备
41  static inline void
42  __rb_rotate_set_parents(struct rb_node *old, struct rb_node *new,
43          struct rb_root *root, int color)
44  {
45      struct rb_node *parent = rb_parent(old);
46      new->__rb_parent_color = old->__rb_parent_color;
47      rb_set_parent_color(old, new, color);
48      __rb_change_child(old, new, parent, root);
49  }
```

linux 红黑树的插入：具体调用为 rb_insert_color, 然后再调用内部函数 __rb_insert, 这里代码太长不在此处展示，具体代码可见 rbtree.c 函数。需要注意的是，linux 的插入函数处理的是节点已经添加在树上后的平衡过程。root 指向红黑树的根节点，但是根节点的 parent 指向 NULL。

linux 红黑树的删除：具体调用 rb_erase 函数，内容如下：

```
1  void rb_erase(struct rb_node *node, struct rb_root *root)
2  {
3      struct rb_node *rebalance;
4      rebalance = __rb_erase_augmented(node, root, &dummy_callbacks);
5      if (rebalance)
```

```
6          _____rb_erase_color(rebalance, root, dummy_rotate);
7  }
8  EXPORT_SYMBOL(rb_erase);
```

___rb_erase_augmented 用于直接删除节点并判断是否破坏了红黑树的性质，在 rbtree_augmented.c 中定义，_____rb_erase_color 用于修复被破坏的红黑树，其代码 rbtree.c

除此之外，linux 红黑树的实现还有替换 rb_replace_node，寻找前驱后继等函数，这些红黑树和 AVL 树没有区别，本实验不需要移植，在这里不再讨论。

## b)  Windows 中的 AVL 树

Windows 中，虚拟内存管理的 VAD 是以 AVL 树的形式管理的。

Windows 的虚拟内存管理利用 MM_AVL_TABLE 型变量, 其定义在 base/ntos/inc/ps.h,

```
1  typedef struct _MM_AVL_TABLE {
2      MMADDRESS_NODE   BalancedRoot;
3      ULONG_PTR DepthOfTree : 5;
4      ULONG_PTR Unused : 3;
5  #if defined (_WIN64)
6      ULONG_PTR NumberGenericTableElements : 56;
7  #else
8      ULONG_PTR NumberGenericTableElements : 24;
9  #endif
10     PVOID NodeHint;
11     PVOID NodeFreeHint;
12 } MM_AVL_TABLE, *PMM_AVL_TABLE;
```

其中的 BalancedRoot 保存了 AVL 树的根节点信息。NumberGenericTableElements 需要在插入和删除节点的时候进行修改。

通过相关资料查询以及对具体源代码的分析可知 BalancedRoot 的 RightChild 指向根节点，而根节点的 Parent 指向 BalancedRoot，这里和 linux 有较大的区别。

MMADDRESS_NODE 的定义在在 base/ntos/mm/mi.h 里:

```
1  typedef struct _MMADDRESS_NODE {
2      union {
3          LONG_PTR Balance : 2;
4          struct _MMADDRESS_NODE *Parent;
5      } u1;
6      struct _MMADDRESS_NODE *LeftChild;
7      struct _MMADDRESS_NODE *RightChild;
8      ULONG_PTR StartingVpn;
9      ULONG_PTR EndingVpn;
```

```
10  } MMADDRESS_NODE, *PMMADDRESS_NODE;
```

和 linux 类似的是，Windows 也是利用存储的父地址的低两位存储节点的性质（平衡因子），不同的是，Windows 的实现方式是联合体。具体实践中发现，该联合体读取时是读出全部 4 字节，但写入时 Balance 和 Parent 需要分开写；红黑树不需要平衡因子，这里利用 Balance 变量存储颜色。后两个数是存储的页信息，和本实验无关。

Windows 中维护 AVL 树的函数主要有：

- 插入 MiInsertNode

- 平衡 MiRebalanceNode 和 MiPromoteNode

- 删除 MiRemoveNode 等等。

对外的接口有 MiInsertNode 和 MiRemoveNode，这也是我们在实验过程中需要保留并重新实现的，它们出现在 /base/ntos/mm/addrsup.c 中。

## 3 实验过程

为了使用 linux 中的函数，先定义以下基础操作

```
1
2  #define rb_black 0
3  #define rb_red 1
4
5  PMMADDRESS_NODE rb_parent(PMMADDRESS_NODE node)
6  {
7      node =SANITIZE_PARENT_NODE(SANITIZE_PARENT_NODE(node)->u1.Parent)
          ;
8      return node;
9  }
10 int rb_color(PMMADDRESS_NODE node)
11 {
12     return node->u1.Balance;
13 }
14 int rb_is_red(PMMADDRESS_NODE node)
15 {
16     return node->u1.Balance==rb_red;
17 }
18 int rb_is_black(PMMADDRESS_NODE node)
19 {
20     return node->u1.Balance==rb_black;
21 }
```

```
22  void rb_set_black(PMMADDRESS_NODE node)
23  {
24      node->u1.Balance = rb_black;
25  }
26  void rb_set_red(PMMADDRESS_NODE node)
27  {
28      node->u1.Balance = rb_red;
29  }
30  void rb_set_parent( PMMADDRESS_NODE rb,PMMADDRESS_NODE p)
31  {
32      rb->u1.Parent =(PMMADDRESS_NODE)(((ULONG_PTR)(p)) + ((ULONG_PTR)(
            rb->u1.Balance)));
33  }
34
35  void rb_set_color( PMMADDRESS_NODE rb, int color)
36  {
37      rb->u1.Balance = color;
38  }
39
40  void rb_set_parent_color(PMMADDRESS_NODE rb, PMMADDRESS_NODE p,int
        color)
41  {
42      rb->u1.Parent =p;
43      rb->u1.Balance=color;
44  }
45
46  void __rb_change_child(PMMADDRESS_NODE old, PMMADDRESS_NODE newer,
47      PMMADDRESS_NODE parent, PMM_AVL_TABLE root)
48  {
49      if (parent!=&root->BalancedRoot) {
50          if (parent->LeftChild== old){
51              parent->LeftChild = newer;
52          }
53          else
54              parent->RightChild = newer;
55      } else
56          root->BalancedRoot.RightChild = newer;
57  }
58
```

```
59  void __rb_rotate_set_parents(PMMADDRESS_NODE old, PMMADDRESS_NODE
        newer,PMM_AVL_TABLE root, int color)
60  {
61      PMMADDRESS_NODE parent = rb_parent(old);
62      rb_set_color(newer,rb_color(old));
63      rb_set_parent(newer,rb_parent(old));
64      rb_set_parent_color(old, newer, color);
65      __rb_change_child(old, newer, parent, root);
66  }
```

其中 SANITIZE_PARENT_NODE 的定义如下 (在 ntos/inc/ps.h 中)

```
1  #define SANITIZE_PARENT_NODE(Parent) ((PMMADDRESS_NODE)(((ULONG_PTR)
        (Parent)) & ~0x3))
```

## a) 插入节点

插入节点的函数在 base/ntos/mm/addrsup.c 中，接口为

```
1  VOID
2  FASTCALL
3  MiInsertNode(
4  IN PMMADDRESS_NODE NodeToInsert,
5  IN PMM_AVL_TABLE Table
6  )
```

与 linux 不同的是，windows 在这一过程中同时进行了插入和平衡的操作，linux 只有平衡，即 MiInsertNode 中插入的一部分应该保留并稍作修改，代码如下：

```
1
2          {
3              PMMADDRESS_NODE NodeOrParent;
4              PMMADDRESS_NODE parent,gparent,tmp;
5              TABLE_SEARCH_RESULT SearchResult;
6              //插入函数
7              SearchResult = MiFindNodeOrParent (Table,
8              NodeToInsert->StartingVpn,
9              &NodeOrParent);
10
11             NodeToInsert->LeftChild = NULL;
12             NodeToInsert->RightChild = NULL;
13
```

```
14              Table->NumberGenericTableElements += 1;

15

16              //
17              // Insert the newer node in the tree.
18              //

19

20              if (SearchResult == TableEmptyTree)
21              {
22                  Table->BalancedRoot.RightChild = NodeToInsert;
23                  rb_set_parent(NodeToInsert,&Table->BalancedRoot);
24              }
25              else
26              {
27                  if (SearchResult == TableInsertAsLeft)
28                  {
29                      NodeOrParent->LeftChild = NodeToInsert;
30                  }
31                  else
32                  {
33                      NodeOrParent->RightChild = NodeToInsert;
34                  }
35                  rb_set_parent(NodeToInsert,NodeOrParent);
36              }
37          rb_set_red(NodeToInsert);
38          parent=rb_parent(NodeToInsert);
39          //平衡部分，直接由linux代码修改，需要注意根节点的性质
40          while(1)  {
41              if (parent==&Table->BalancedRoot) {
42              Table->BalancedRoot.RightChild = NodeToInsert;
43              rb_set_parent_color(NodeToInsert,&Table->BalancedRoot
                    ,rb_black);
44              break;
45          } else if (rb_is_black(parent))
46          break;
47          gparent = rb_parent(parent);
48          tmp = gparent->RightChild;

49

50          if (parent != tmp) {     /* parent == gparent->rb_left */
51              if (tmp && rb_is_red(tmp)) {
```

```
52              /*
53               * Case 1 − color flips
54               *
55               *        G            g
56               *       / \          / \
57               *      p   u  −−>   P   U
58               *     /            /
59               *    n            n
60               *
61               * However, since g's parent might be red, and
62               * 4) does not allow this, we need to recurse
63               * at g.
64               */
65              rb_set_parent_color(tmp, gparent, rb_black);
66              rb_set_parent_color(parent, gparent, rb_black);
67              NodeToInsert = gparent;
68              parent = rb_parent(NodeToInsert);
69              rb_set_parent_color(NodeToInsert, parent, rb_red);
70              continue;
71          }

73          tmp = parent−>RightChild;
74          if (NodeToInsert == tmp) {
75              /*
76               * Case 2 − left rotate at parent
77               *
78               *       G            G
79               *      / \          / \
80               *     p   U  −−>   n   U
81               *      \          /
82               *       n        p
83               *
84               * This still leaves us in violation of 4), the
85               * continuation into Case 3 will fix that.
86               */
87              parent−>RightChild = tmp = NodeToInsert−>LeftChild;
88              NodeToInsert−>LeftChild = parent;
89              if (tmp)
90              rb_set_parent_color(tmp, parent,rb_black);
```

```
91              rb_set_parent_color(parent, NodeToInsert, rb_red);
92              parent = NodeToInsert;
93              tmp = NodeToInsert->RightChild;
94          }
95
96          /*
97           * Case 3 - right rotate at gparent
98           *
99           *        G           P
100          *       / \         / \
101          *      p   U  -->  n   g
102          *     /                 \
103          *    n                   U
104          */
105         gparent->LeftChild= tmp;   /* == parent->rb_right */
106         parent->RightChild= gparent;
107         if (tmp)
108         rb_set_parent_color(tmp, gparent, rb_black);
109         __rb_rotate_set_parents(gparent, parent, Table, rb_red);
110         break;
111     } else {
112         tmp = gparent->LeftChild;
113         if (tmp && rb_is_red(tmp)) {
114         /* Case 1 - color flips */
115         rb_set_parent_color(tmp, gparent, rb_black);
116         rb_set_parent_color(parent, gparent, rb_black);
117         NodeToInsert = gparent;
118         parent = rb_parent(NodeToInsert);
119         rb_set_parent_color(NodeToInsert, parent, rb_red);
120         continue;
121     }
122
123     tmp = parent->LeftChild;
124     if (NodeToInsert == tmp) {
125         /* Case 2 - right rotate at parent */
126         parent->LeftChild = tmp = NodeToInsert->RightChild;
127         NodeToInsert->RightChild = parent;
128         if (tmp)
129         rb_set_parent_color(tmp, parent,
```

```
130                rb_black);
131                rb_set_parent_color(parent, NodeToInsert, rb_red);
132                parent = NodeToInsert;
133                tmp = NodeToInsert->LeftChild;
134            }
135
136            /* Case 3 - left rotate at gparent */
137            gparent->RightChild = tmp;   /* == parent->rb_left */
138            parent->LeftChild = gparent;
139            if (tmp)
140            rb_set_parent_color(tmp, gparent, rb_black);
141            __rb_rotate_set_parents(gparent, parent, Table, rb_red);
142            break;
143        }
144    }
145 return;
146 }
```

有趣的是，linux 中给 ___rb_insert 传入 void (*augment_rotate) 这一函数指针，但是该函数是空的（应该是给用户调用提供类似重载的特性），所以直接将相关语句删除即可。

## b)  删除节点

删除节点的函数也出现在 base/ntos/mm/addrsup.c 中，其接口为：

```
1 VOID
2 FASTCALL
3 MiRemoveNode(
4 IN PMMADDRESS_NODE NodeToDelete,
5 IN PMM_AVL_TABLE Table
6 )
```

使用 linux 代码前需要先移植 ___rb_erase_augmente 和 _____rb_erase_color; 代码如下

```
1
2 static PMMADDRESS_NODE ___rb_erase_augmented(PMMADDRESS_NODE node,
       PMM_AVL_TABLE root)
3 {
4     PMMADDRESS_NODE   child = node->RightChild, tmp = node->LeftChild;
5     PMMADDRESS_NODE   parent, rebalance;
6     PMMADDRESS_NODE pc;
7
```

```
 8     if (!tmp) {
 9         /*
10          * Case 1: node to erase has no more than 1 child (easy!)
11          *
12          * Note that if there is one child it must be red due to 5)
13          * and node must be black due to 4). We adjust colors locally
14          * so as to bypass ___rb_erase_color() later on.
15          */
16         pc = node->u1.Parent;
17         parent = rb_parent(node);
18         ___rb_change_child(node, child, parent, root);
19         if (child) {
20             child->u1.Parent = pc;
21             child->u1.Balance=(((ULONG_PTR)(pc)) & 0x1);
22             rebalance = &root->BalancedRoot;
23         } else
24             rebalance = (((ULONG_PTR)(pc)) & 0x1)==0 ? parent : &root
                    ->BalancedRoot;
25     } else if (!child) {
26         /* Still case 1, but this time the child is node->rb_left */
27         tmp->u1.Parent =pc= node->u1.Parent;
28         tmp->u1.Balance=node->u1.Balance;
29         parent = rb_parent(node);
30         ___rb_change_child(node, tmp, parent, root);
31         rebalance = &root->BalancedRoot;
32     } else {
33         PMMADDRESS_NODE successor = child, child2;
34         tmp = child->LeftChild;
35         if (!tmp) {
36             /*
37              * Case 2: node's successor is its right child
38              *
39              *    (n)          (s)
40              *    / \          / \
41              *  (x) (s)  ->  (x) (c)
42              *        \
43              *        (c)
44              */
45             parent = successor;
```

```
46            child2 = successor->RightChild;
47        } else {
48            /*
49             * Case 3: node's successor is leftmost under
50             * node's right child subtree
51             *
52             *     (n)          (s)
53             *     / \          / \
54             *   (x) (y)   ->  (x) (y)
55             *       /              /
56             *     (p)          (p)
57             *     /              /
58             *   (s)          (c)
59             *     \
60             *     (c)
61             */
62            do {
63                parent = successor;
64                successor = tmp;
65                tmp = tmp->LeftChild;
66            } while (tmp);
67            parent->LeftChild = child2 = successor->RightChild;
68            successor->RightChild = child;
69            rb_set_parent(child, successor);
70        }
71
72        successor->LeftChild = tmp = node->LeftChild;
73        rb_set_parent(tmp, successor);
74        pc = node->u1.Parent;
75        tmp = SANITIZE_PARENT_NODE(pc);
76        __rb_change_child(node, successor, tmp, root);
77        if (child2) {
78            successor->u1.Parent = pc;
79            successor->u1.Balance = (((ULONG_PTR)(pc)) & 0x1);
80            rb_set_parent_color(child2, parent, rb_black);
81            rebalance = &root->BalancedRoot;
82        } else {
83            PMMADDRESS_NODE pc2 = successor->u1.Parent;
84            successor->u1.Parent = pc;
```

```
85              successor->u1.Balance = (((ULONG_PTR)(pc)) & 0x1);
86              rebalance = (((ULONG_PTR)(pc2)) & 0x1)==0 ? parent : &
                    root->BalancedRoot;
87          }
88          tmp = successor;
89      }
90      return rebalance;
91  }
92
93  static void
94      _____rb_erase_color(PMMADDRESS_NODE   parent, PMM_AVL_TABLE root)
95  {
96      PMMADDRESS_NODE node =NULL, sibling, tmp1=NULL, tmp2=NULL;
97
98      while (1) {
99          /*
100         * Loop invariants:
101         * - node is black (or NULL on first iteration)
102         * - node is not the root (parent is not NULL)
103         * - All leaf paths going through parent and node have a
104         *   black node count that is 1 lower than other leaf paths.
105         */
106         sibling = parent->RightChild;
107         if (node != sibling) {   /* node == parent->rb_left */
108             if (rb_is_red(sibling)) {
109                 /*
110                 * Case 1 - left rotate at parent
111                 *
112                 *     P               S
113                 *    / \             / \
114                 *   N   s    -->    p   Sr
115                 *      / \         / \
116                 *     Sl  Sr      N   Sl
117                 */
118                 parent->RightChild = tmp1 = sibling->LeftChild;
119                 sibling->LeftChild = parent;
120                 if(tmp1)
121                     rb_set_parent_color(tmp1, parent, rb_black);
122                 __rb_rotate_set_parents(parent, sibling, root,
```

```
123                          rb_red);
124                  sibling = tmp1;
125              }
126          tmp1 = sibling->RightChild;
127          if (!tmp1 || rb_is_black(tmp1)) {
128                  tmp2 = sibling->LeftChild;
129              if (!tmp2 || rb_is_black(tmp2)) {
130                  /*
131                   * Case 2 - sibling color flip
132                   * (p could be either color here)
133                   *
134                   *    (p)           (p)
135                   *    / \           / \
136                   *   N   S   -->   N   s
137                   *      / \           / \
138                   *     Sl  Sr        Sl  Sr
139                   *
140                   * This leaves us violating 5) which
141                   * can be fixed by flipping p to black
142                   * if it was red, or by recursing at p.
143                   * p is red when coming from Case 1.
144                   */
145                      rb_set_parent_color(sibling, parent,
146                      rb_red);
147                  if (rb_is_red(parent))
148                      rb_set_black(parent);
149                  else {
150                      node = parent;
151                      parent = rb_parent(node);
152                      if (parent!=&root->BalancedRoot)
153                          continue;
154                  }
155                  break;
156              }
157              /*
158               * Case 3 - right rotate at sibling
159               * (p could be either color here)
160               *
161               *    (p)           (p)
```

```
162              *    / \            / \
163              *   N    S     —>  N    Sl
164              *        / \             \
165              *       sl   Sr            s
166              *                           \
167              *                            Sr
168              */
169              sibling->LeftChild = tmp1 = tmp2->RightChild;
170              tmp2->RightChild = sibling;
171              parent->RightChild = tmp2;
172              if (tmp1)
173                  rb_set_parent_color(tmp1, sibling,
174                  rb_black);
175              tmp1 = sibling;
176              sibling = tmp2;
177          }
178          /*
179          * Case 4 - left rotate at parent + color flips
180          * (p and sl could be either color here.
181          *  After rotation, p becomes black, s acquires
182          *  p's color, and sl keeps its color)
183          *
184          *      (p)            (s)
185          *      / \            / \
186          *    N    S     —>   P    Sr
187          *        / \        / \
188          *      (sl) sr     N  (sl)
189          */
190          parent->RightChild = tmp2 = sibling->LeftChild;
191          sibling->LeftChild = parent;
192          rb_set_parent_color(tmp1, sibling, rb_black);
193          if (tmp2)
194              rb_set_parent(tmp2, parent);
195          ___rb_rotate_set_parents(parent, sibling, root,
196              rb_black);
197          break;
198      } else {
199          sibling = parent->LeftChild;
200          if (rb_is_red(sibling)) {
```

```
201                    /* Case 1 - right rotate at parent */
202                    parent->LeftChild = tmp1 = sibling->RightChild;
203                    sibling->RightChild = parent;
204                    rb_set_parent_color(tmp1, parent, rb_black);
205                    __rb_rotate_set_parents(parent, sibling, root,
206                        rb_red);
207                    sibling = tmp1;
208                }
209                tmp1 = sibling->LeftChild;
210                if (!tmp1 || rb_is_black(tmp1)) {
211                    tmp2 = sibling->RightChild;
212                    if (!tmp2 || rb_is_black(tmp2)) {
213                        /* Case 2 - sibling color flip */
214                        rb_set_parent_color(sibling, parent,
215                            rb_red);
216                        if (rb_is_red(parent))
217                            rb_set_black(parent);
218                        else {
219                            node = parent;
220                            parent = rb_parent(node);
221                            if (parent!=&root->BalancedRoot)
222                                continue;
223                        }
224                        break;
225                    }
226                    /* Case 3 - right rotate at sibling */
227                    sibling->RightChild = tmp1 = tmp2->LeftChild;
228                    tmp2->LeftChild = sibling;
229                    parent->LeftChild = tmp2;
230                    if (tmp1)
231                        rb_set_parent_color(tmp1, sibling,
232                        rb_black);
233                    tmp1 = sibling;
234                    sibling = tmp2;
235                }
236                /* Case 4 - left rotate at parent + color flips */
237                parent->LeftChild = tmp2 = sibling->RightChild;
238                sibling->RightChild = parent;
239                rb_set_parent_color(tmp1, sibling, rb_black);
```

```
240        if (tmp2)
241            rb_set_parent(tmp2, parent);
242        __rb_rotate_set_parents(parent, sibling, root,
243            rb_black);
244        break;
245        }
246    }
247 }
```

值得注意的是根节点的 Parent 值为 &Table->BalancedRoot, 以及记录颜色和父节点信息的联合体的特性。

之后再修改 MiRemoveNode 函数为,

```
1 {
2    PMMADDRESS_NODE rebalance;
3    Table->NumberGenericTableElements -= 1;
4    rebalance = __rb_erase_augmented(NodeToDelete, Table);
5    if (rebalance!=&Table->BalancedRoot&&rebalance)
6        _____rb_erase_color(rebalance, Table);
7 }
```

即可完成移植

# 4  实验结果

到此为止，代码移植工作已经完成，编译好后 (nmake -nologo x86=) 复制到虚拟机（Windows Server 2003 SP2）的 C:/Windows/system32 文件夹，设置好主机的 WinDbg[1]调试器和虚拟机的启动引导信息[2]，然互选择 Debug 引导启动, 如图1。

一切正常则能够正常进入系统，并可以正常操作，如图2。通过 WinDbg 可以发现，在系统的初始或与平常使用中，虚拟内存管理是一直使用的，在对该红黑树在不断的插入与删除节点，所以系统可以稳定运行就意味着没有 bug。

# 5  实验感想

本实验一开始思考比较简单，因为无论是 WRK 还是 Linux 内核都有着大量的资料，而且本实验的目的是一致，只需要利用原有的接口进行一定的修改即可。

但是随着任务的深入，很多问题并不如想象的顺利：linux 在 3.x 版本起对红黑树部分进行了大改，原有文档过于陈旧；系统接口难以的定位；MS 的 symbol 服务器关机导致 WinDbg 难以进行复杂调试；

---

[1]见 debug.bat 文件, 需先导入 debug.WEW

[2]boot.in 增加信息 multi(0)disk(0)rdisk(0)partition(1) \WINDOWS="Debug" /kernel=wrkx86.exe /hal=halmacpi.dll /debug /debugport=com1 /baudrate=115200
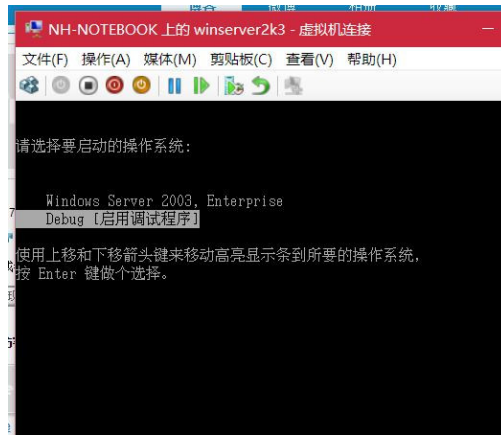
图 1: boot 选项



图 2: 正常工作

两者数据结构差异而导致的算法细节区别较大；联合体这一少见数据结构的陷阱。

这一个个问题很令我烦恼，也消耗了大量的时间在其中去解决，但我从中学到了很多。我学会了更好的利用 ctag 去阅读代码，学会了更好的从大量资料中去寻找信息，学会了将大型项目中的函数及其依赖剥离出来进行测试[3]，了解了 Linux 内核和 Windows 内核迥异的代码风格以及他们在虚拟内存管上的异同，对操作系统有了更深一步的理解。

虽然本次实验所涉及的 VAD 修改只是操作系统的一小部分，它仅仅完成是内存的组织、分配、调度和回收，但这一过程除去了操作系统的神秘感，使我能更好的去探索这两份代码中的秘密，也为以后的学习工作带来很好的基础。

# 6　实验代码

系统原有代码见 origin 文件夹，修改后的部分都存储在 addrsup.c 文件中, 使用时直接用 addrsup.c 替换 wrk 中的同名文件编译即可。

```
1   /*++
2
3   Copyright (c) Microsoft Corporation. All rights reserved.
4
5   You may only use this code if you agree to the terms of the Windows
        Research Kernel Source Code License agreement (see License.txt).
6   If you do not agree to the terms, do not use the code.
7
8
9   Module Name:
10
11      addrsup.c
12
13  Abstract:
14
15      This module implements a new version of the generic table package
16      based on balanced binary trees (later named AVL), as described in
17      Knuth, "The Art of Computer Programming, Volume 3, Sorting and
            Searching",
18      and refers directly to algorithms as they are presented in the
            second
19      edition Copyrighted in 1973.
20
21      Used rtl\avltable.c as a starting point, adding the following:
22
```

---

[3]见测试文件中的测试项目

```
23        - Use less memory for structures as these are nonpaged & heavily
             used .
24        - Caller allocates the pool to reduce mutex hold times .
25        - Various VAD-specific customizations/optimizations .
26        - Hints .

27
28   Environment :

29
30        Kernel mode only , working set mutex held , APCs disabled .

31
32   --*/

33
34   #include "mi.h"

35
36   #define rb_black 0
37   #define rb_red 1

38
39   PMMADDRESS_NODE rb_parent (PMMADDRESS_NODE node )
40   {
41        node =SANITIZE_PARENT_NODE(SANITIZE_PARENT_NODE( node)->u1 . Parent )
             ;
42        return node ;
43   }
44   int rb_color (PMMADDRESS_NODE node )
45   {
46        return node->u1 . Balance ;
47   }
48   int rb_is_red (PMMADDRESS_NODE node )
49   {
50        return node->u1 . Balance==rb_red ;
51   }
52   int rb_is_black (PMMADDRESS_NODE node )
53   {
54        return node->u1 . Balance==rb_black ;
55   }
56   void rb_set_black (PMMADDRESS_NODE node )
57   {
58        node->u1 . Balance = rb_black ;
59   }
```

```
60  void rb_set_red(PMMADDRESS_NODE node)
61  {
62      node->u1.Balance = rb_red;
63  }
64  void rb_set_parent( PMMADDRESS_NODE rb ,PMMADDRESS_NODE p)
65  {
66      rb->u1.Parent =(PMMADDRESS_NODE)(((ULONG_PTR)(p)) + ((ULONG_PTR)(
            rb->u1.Balance)));
67  }
68
69  void rb_set_color( PMMADDRESS_NODE rb, int color)
70  {
71      rb->u1.Balance = color;
72  }
73
74  void rb_set_parent_color(PMMADDRESS_NODE rb, PMMADDRESS_NODE p,int
        color)
75  {
76      rb->u1.Parent =p;
77      rb->u1.Balance=color;
78  }
79
80  void
81      __rb_change_child(PMMADDRESS_NODE old , PMMADDRESS_NODE newer ,
82      PMMADDRESS_NODE parent , PMM_AVL_TABLE root)
83  {
84      if (parent!=&root->BalancedRoot) {
85          if (parent->LeftChild== old){
86              parent->LeftChild = newer;
87          }
88          else
89              parent->RightChild = newer;
90      } else
91          root->BalancedRoot.RightChild = newer;
92  }
93
94  void __rb_rotate_set_parents(PMMADDRESS_NODE old , PMMADDRESS_NODE
        newer ,
95                                  PMM_AVL_TABLE root , int color)
```

```
96  {
97       PMMADDRESS_NODE parent = rb_parent(old);
98       rb_set_color(newer,rb_color(old));
99       rb_set_parent(newer,rb_parent(old));
100      rb_set_parent_color(old, newer, color);
101      __rb_change_child(old, newer, parent, root);
102  }
103
104  static PMMADDRESS_NODE __rb_erase_augmented(PMMADDRESS_NODE node,
         PMM_AVL_TABLE root)
105  {
106      PMMADDRESS_NODE  child = node->RightChild, tmp = node->LeftChild;
107      PMMADDRESS_NODE  parent, rebalance;
108      PMMADDRESS_NODE pc;
109
110      if (!tmp) {
111          /*
112          * Case 1: node to erase has no more than 1 child (easy!)
113          *
114          * Note that if there is one child it must be red due to 5)
115          * and node must be black due to 4). We adjust colors locally
116          * so as to bypass __rb_erase_color() later on.
117          */
118          pc = node->u1.Parent;
119          parent = rb_parent(node);
120          __rb_change_child(node, child, parent, root);
121          if (child) {
122              child->u1.Parent = pc;
123              child->u1.Balance=(((ULONG_PTR)(pc)) & 0x1);
124              rebalance = &root->BalancedRoot;
125          } else
126              rebalance = (((ULONG_PTR)(pc)) & 0x1)==0 ? parent : &root
                    ->BalancedRoot;
127      } else if (!child) {
128          /* Still case 1, but this time the child is node->rb_left */
129          tmp->u1.Parent =pc= node->u1.Parent;
130          tmp->u1.Balance=node->u1.Balance;
131          parent = rb_parent(node);
132          __rb_change_child(node, tmp, parent, root);
```

```
133             rebalance = &root->BalancedRoot;
134     } else {
135         PMMADDRESS_NODE successor = child, child2;
136         tmp = child->LeftChild;
137         if (!tmp) {
138             /*
139             * Case 2: node's successor is its right child
140             *
141             *     (n)          (s)
142             *     / \          / \
143             *   (x) (s)  ->  (x) (c)
144             *         \
145             *         (c)
146             */
147             parent = successor;
148             child2 = successor->RightChild;
149         } else {
150             /*
151             * Case 3: node's successor is leftmost under
152             * node's right child subtree
153             *
154             *     (n)          (s)
155             *     / \          / \
156             *   (x) (y)  ->  (x) (y)
157             *       /            /
158             *     (p)          (p)
159             *     /            /
160             *   (s)          (c)
161             *     \
162             *     (c)
163             */
164             do {
165                 parent = successor;
166                 successor = tmp;
167                 tmp = tmp->LeftChild;
168             } while (tmp);
169             parent->LeftChild = child2 = successor->RightChild;
170             successor->RightChild = child;
171             rb_set_parent(child, successor);
```

```
172                 }
173
174                 successor->LeftChild = tmp = node->LeftChild;
175                 rb_set_parent(tmp, successor);
176                 pc = node->u1.Parent;
177                 tmp = SANITIZE_PARENT_NODE(pc);
178                 __rb_change_child(node, successor, tmp, root);
179                 if (child2) {
180                     successor->u1.Parent = pc;
181                     successor->u1.Balance = (((ULONG_PTR)(pc)) & 0x1);
182                     rb_set_parent_color(child2, parent, rb_black);
183                     rebalance = &root->BalancedRoot;
184                 } else {
185                     PMMADDRESS_NODE pc2 = successor->u1.Parent;
186                     successor->u1.Parent = pc;
187                     successor->u1.Balance = (((ULONG_PTR)(pc)) & 0x1);
188                     rebalance = (((ULONG_PTR)(pc2)) & 0x1)==0 ? parent : &
                            root->BalancedRoot;
189                 }
190                 tmp = successor;
191         }
192     return rebalance;
193 }
194
195 static void
196     _____rb_erase_color(PMMADDRESS_NODE  parent, PMM_AVL_TABLE root)
197 {
198     PMMADDRESS_NODE node =NULL, sibling, tmp1=NULL, tmp2=NULL;
199
200     while (1) {
201         /*
202         * Loop invariants:
203         * - node is black (or NULL on first iteration)
204         * - node is not the root (parent is not NULL)
205         * - All leaf paths going through parent and node have a
206         *   black node count that is 1 lower than other leaf paths.
207         */
208         sibling = parent->RightChild;
209         if (node != sibling) {  /* node == parent->rb_left */
```

```
210                if (rb_is_red(sibling)) {
211                    /*
212                    * Case 1 - left rotate at parent
213                    *
214                    *     P                S
215                    *    / \              / \
216                    *   N   s    -->     p   Sr
217                    *      / \          / \
218                    *     Sl  Sr       N   Sl
219                    */
220                parent->RightChild = tmp1 = sibling->LeftChild;
221                sibling->LeftChild = parent;
222                if(tmp1)
223                    rb_set_parent_color(tmp1, parent, rb_black);
224                __rb_rotate_set_parents(parent, sibling, root,
225                    rb_red);
226                sibling = tmp1;
227            }
228            tmp1 = sibling->RightChild;
229            if (!tmp1 || rb_is_black(tmp1)) {
230                    tmp2 = sibling->LeftChild;
231                if (!tmp2 || rb_is_black(tmp2)) {
232                    /*
233                    * Case 2 - sibling color flip
234                    * (p could be either color here)
235                    *
236                    *     (p)              (p)
237                    *     / \              / \
238                    *    N   S    -->     N   s
239                    *       / \              / \
240                    *      Sl  Sr           Sl  Sr
241                    *
242                    * This leaves us violating 5) which
243                    * can be fixed by flipping p to black
244                    * if it was red, or by recursing at p.
245                    * p is red when coming from Case 1.
246                    */
247                    rb_set_parent_color(sibling, parent,
248                        rb_red);
```

```
249                    if (rb_is_red(parent))
250                        rb_set_black(parent);
251                    else {
252                        node = parent;
253                        parent = rb_parent(node);
254                        if (parent!=&root->BalancedRoot)
255                            continue;
256                    }
257                    break;
258                }
259                /*
260                 * Case 3 - right rotate at sibling
261                 * (p could be either color here)
262                 *
263                 *    (p)           (p)
264                 *    / \           / \
265                 *   N   S    -->  N   Sl
266                 *      / \             \
267                 *     sl  Sr            s
268                 *                        \
269                 *                         Sr
270                 */
271                sibling->LeftChild = tmp1 = tmp2->RightChild;
272                tmp2->RightChild = sibling;
273                parent->RightChild = tmp2;
274                if (tmp1)
275                    rb_set_parent_color(tmp1, sibling ,
276                    rb_black);
277                tmp1 = sibling;
278                sibling = tmp2;
279            }
280            /*
281             * Case 4 - left rotate at parent + color flips
282             * (p and sl could be either color here.
283             *  After rotation, p becomes black, s acquires
284             *  p's color, and sl keeps its color)
285             *
286             *      (p)              (s)
287             *      / \              / \
```

```
288              *      N    S      —>     P    Sr
289              *           / \            / \
290              *          (sl) sr       N   (sl)
291              */
292              parent->RightChild = tmp2 = sibling->LeftChild;
293              sibling->LeftChild = parent;
294              rb_set_parent_color(tmp1, sibling, rb_black);
295              if (tmp2)
296                  rb_set_parent(tmp2, parent);
297              __rb_rotate_set_parents(parent, sibling, root,
298                  rb_black);
299              break;
300          } else {
301              sibling = parent->LeftChild;
302              if (rb_is_red(sibling)) {
303                  /* Case 1 – right rotate at parent */
304                  parent->LeftChild = tmp1 = sibling->RightChild;
305                  sibling->RightChild = parent;
306                  rb_set_parent_color(tmp1, parent, rb_black);
307                  __rb_rotate_set_parents(parent, sibling, root,
308                      rb_red);
309                  sibling = tmp1;
310              }
311              tmp1 = sibling->LeftChild;
312              if (!tmp1 || rb_is_black(tmp1)) {
313                  tmp2 = sibling->RightChild;
314                  if (!tmp2 || rb_is_black(tmp2)) {
315                      /* Case 2 – sibling color flip */
316                      rb_set_parent_color(sibling, parent,
317                          rb_red);
318                      if (rb_is_red(parent))
319                          rb_set_black(parent);
320                      else {
321                          node = parent;
322                          parent = rb_parent(node);
323                          if (parent!=&root->BalancedRoot)
324                              continue;
325                      }
326                      break;
```

```
327                     }
328                     /* Case 3 - right rotate at sibling */
329                     sibling->RightChild = tmp1 = tmp2->LeftChild;
330                     tmp2->LeftChild = sibling;
331                     parent->LeftChild = tmp2;
332                     if (tmp1)
333                         rb_set_parent_color(tmp1, sibling,
334                         rb_black);
335                     tmp1 = sibling;
336                     sibling = tmp2;
337                 }
338                 /* Case 4 - left rotate at parent + color flips */
339                 parent->LeftChild = tmp2 = sibling->RightChild;
340                 sibling->RightChild = parent;
341                 rb_set_parent_color(tmp1, sibling, rb_black);
342                 if (tmp2)
343                     rb_set_parent(tmp2, parent);
344                 __rb_rotate_set_parents(parent, sibling, root,
345                     rb_black);
346                 break;
347             }
348         }
349 }
350
351 #if !defined (_USERMODE)
352 #define PRINT
353 #define COUNT_BALANCE_MAX(a)
354 #else
355 extern MM_AVL_TABLE MmSectionBasedRoot;
356 #endif
357
358 #if (_MSC_VER >= 800)
359 #pragma warning(disable:4010)         // Allow pretty pictures without
         the noise
360 #endif
361
362 TABLE_SEARCH_RESULT
363 MiFindNodeOrParent (
364     IN PMM_AVL_TABLE Table,
```

```
365        IN ULONG_PTR StartingVpn ,
366        OUT PMMADDRESS_NODE *NodeOrParent
367        ) ;
368
369  VOID
370  MiPromoteNode (
371        IN PMMADDRESS_NODE C
372        ) ;
373
374  ULONG
375  MiRebalanceNode (
376        IN PMMADDRESS_NODE S
377        ) ;
378
379  PMMADDRESS_NODE
380  MiRealSuccessor (
381        IN PMMADDRESS_NODE Links
382        ) ;
383
384  PMMADDRESS_NODE
385  MiRealPredecessor (
386        IN PMMADDRESS_NODE Links
387        ) ;
388
389  VOID
390  MiInitializeVadTableAvl (
391        IN PMM_AVL_TABLE Table
392        ) ;
393
394  PVOID
395  MiEnumerateGenericTableWithoutSplayingAvl (
396        IN PMM_AVL_TABLE Table ,
397        IN PVOID *RestartKey
398        ) ;
399
400  #ifdef ALLOC_PRAGMA
401  #pragma alloc_text (PAGE, MiCheckForConflictingNode )
402  #pragma alloc_text (PAGE, MiRealSuccessor )
403  #pragma alloc_text (PAGE, MiRealPredecessor )
```

```
404  #pragma  alloc_text(PAGE, MiInitializeVadTableAvl)
405  #pragma  alloc_text(PAGE, MiFindEmptyAddressRangeInTree)
406  #pragma  alloc_text(PAGE, MiFindEmptyAddressRangeDownTree)
407  #pragma  alloc_text(PAGE, MiFindEmptyAddressRangeDownBasedTree)
408  #endif
409
410  //
411  //  Various  Rtl  macros  that  reference  Parent  use  private  versions  here
            since
412  //  Parent  is  overloaded  with  Balance.
413  //
414
415  //
416  //   The  macro  function  Parent  takes  as  input  a  pointer  to  a  splay
           link  in  a
417  //   tree  and  returns  a  pointer  to  the  splay  link  of  the  parent  of  the
           input
418  //   node.   If  the  input  node  is  the  root  of  the  tree  the  return  value
           is
419  //   equal  to  the  input  value.
420  //
421  //   PRTL_SPLAY_LINKS
422  //   MiParent  (
423  //        PRTL_SPLAY_LINKS  Links
424  //        );
425  //
426
427  #define  MiParent(Links)  (                        \
428      (PRTL_SPLAY_LINKS)(SANITIZE_PARENT_NODE((Links)->u1.Parent))  \
429      )
430
431  //
432  //   The  macro  function  IsLeftChild  takes  as  input  a  pointer  to  a
           splay  link
433  //   in  a  tree  and  returns  TRUE  if  the  input  node  is  the  left  child  of
           its
434  //   parent,  otherwise  it  returns  FALSE.
435  //
436  //   BOOLEAN
```

```
437  //   MiIsLeftChild (
438  //        PRTL_SPLAY_LINKS Links
439  //        );
440  //

442  #define MiIsLeftChild(Links) (                                    \
443      (RtlLeftChild(MiParent(Links)) == (PRTL_SPLAY_LINKS)(Links)) \
444      )


446  //
447  //   The macro function IsRightChild takes as input a pointer to a
         splay link
448  //  in a tree and returns TRUE if the input node is the right child
         of its
449  //  parent, otherwise it returns FALSE.
450  //
451  //   BOOLEAN
452  //   MiIsRightChild (
453  //        PRTL_SPLAY_LINKS Links
454  //        );
455  //

457  #define MiIsRightChild(Links) (                                    \
458      (RtlRightChild(MiParent(Links)) == (PRTL_SPLAY_LINKS)(Links)) \
459      )



463  #if DBG

465  //
466  // Build a table of the best case efficiency of a balanced binary
         tree,
467  // holding the most possible nodes that can possibly be held in a
         binary
468  // tree with a given number of levels.  The answer is always (2**n) -
          1.
469  //
470  // (Used for debug only.)
```

```
471  //
472
473  ULONG  MiBestCaseFill[33] = {
474          0,              1,              3,              7,
475          0xf,            0x1f,           0x3f,           0x7f,
476          0xff,           0x1ff,          0x3ff,          0x7ff,
477          0xfff,          0x1fff,         0x3fff,         0x7fff,
478          0xffff,         0x1ffff,        0x3ffff,        0x7ffff,
479          0xfffff,        0x1fffff,       0x3fffff,       0x7fffff,
480          0xffffff,       0x1ffffff,      0x3ffffff,      0x7ffffff,
481          0xfffffff,      0x1fffffff,     0x3fffffff,     0x7fffffff,
482          0xffffffff
483  };
484
485  //
486  // Build a table of the worst case efficiency of a balanced binary
         tree,
487  // holding the fewest possible nodes that can possibly be contained
         in a
488  // balanced binary tree with the given number of levels.  After the
         first
489  // two levels, each level n is obviously occupied by a root node,
         plus
490  // one subtree the size of level n−1, and another subtree which is
         the
491  // size of n−2, i.e.:
492  //
493  //      MiWorstCaseFill[n] = 1 + MiWorstCaseFill[n−1] +
         MiWorstCaseFill[n−2]
494  //
495  // The efficiency of a typical balanced binary tree will normally
         fall
496  // between the two extremes, typically closer to the best case.  Note
497  // however that even with the worst case, it only takes 32 compares
         to
498  // find an element in a worst case tree populated with ~3.5M nodes.
499  //
500  // Unbalanced trees and splay trees, on the other hand, can and will
         sometimes
```

```
501  // degenerate to a straight line, requiring on average n/2 compares
         to
502  // find a node.
503  //
504  // A specific case is one where the nodes are inserted in collated
         order.
505  // In this case an unbalanced or a splay tree will generate a
         straight
506  // line, yet the balanced binary tree will always create a perfectly
507  // balanced tree (best-case fill) in this situation.
508  //
509  // (Used for debug only.)
510  //
511
512  ULONG MiWorstCaseFill[33] = {
513          0,              1,              2,              4,
514          7,              12,             20,             33,
515          54,             88,             143,            232,
516          376,            609,            986,            1596,
517          2583,           4180,           6764,           10945,
518          17710,          28656,          46367,          75024,
519          121392,         196417,         317810,         514228,
520          832039,         1346268,        2178308,        3524577,
521          5702886
522  };
523
524  #endif
525
526
527  TABLE_SEARCH_RESULT
528  MiFindNodeOrParent (
529      IN PMM_AVL_TABLE Table,
530      IN ULONG_PTR StartingVpn,
531      OUT PMMADDRESS_NODE *NodeOrParent
532      )
533
534  /*++
535
536  Routine Description:
```

```
537
538        This  routine  is  used  by  all  of  the  routines  of  the  generic
539        table  package  to  locate  the  a  node  in  the  tree.   It  will
540        find  and  return  (via  the  NodeOrParent  parameter)  the  node
541        with  the  given  key,  or  if  that  node  is  not  in  the  tree  it
542        will  return  (via  the  NodeOrParent  parameter)  a  pointer  to
543        the  parent.
544
545   Arguments:
546
547        Table − The  generic  table  to  search  for  the  key.
548
549        StartingVpn − The  starting  virtual  page  number.
550
551        NodeOrParent − Will  be  set  to  point  to  the  node  containing  the
552                          the  key  or  what  should  be  the  parent  of  the  node
553                          if  it  were  in  the  tree.   Note  that  this  will  *NOT*
554                          be  set  if  the  search  result  is  TableEmptyTree.
555
556   Return  Value:
557
558        TABLE_SEARCH_RESULT − TableEmptyTree:  The  tree  was  empty.
               NodeOrParent
559                                                       is  *not*  altered.
560
561                                 TableFoundNode:  A  node  with  the  key  is  in
                                       the  tree.
562                                                   NodeOrParent  points  to  that
                                                         node.
563
564                                 TableInsertAsLeft:  Node  with  key  was  not
                                       found.
565                                                   NodeOrParent  points  to
                                                         what  would
566                                                   be  parent.   The  node
                                                         would  be  the
567                                                   left  child.
568
569                                 TableInsertAsRight:  Node  with  key  was  not
```

```
                                        found.
570                                                 NodeOrParent points to
                                                        what would
571                                                 be parent.  The node
                                                        would be
572                                                 the right child.

574  Environment:

576       Kernel mode.  The PFN lock is held for some of the tables.

578  ––*/

580  {
581  #if DBG
582       ULONG NumberCompares = 0;
583  #endif
584       PMMADDRESS_NODE Child;
585       PMMADDRESS_NODE NodeToExamine;

587       if (Table–>NumberGenericTableElements == 0) {
588           return TableEmptyTree;
589       }

591       NodeToExamine = (PMMADDRESS_NODE) Table–>BalancedRoot.RightChild;

593       do {

595           //
596           // Make sure the depth of tree is correct.
597           //

599           // ASSERT(++NumberCompares <= Table–>DepthOfTree);

601           //
602           // Compare the buffer with the key in the tree element.
603           //

605           if (StartingVpn < NodeToExamine–>StartingVpn) {
```

```
606
607                 Child = NodeToExamine->LeftChild;

608
609                 if (Child != NULL) {
610                     NodeToExamine = Child;
611                 }
612                 else {

613
614                     //
615                     // Node is not in the tree.  Set the output
616                     // parameter to point to what would be its
617                     // parent and return which child it would be.
618                     //

619
620                     *NodeOrParent = NodeToExamine;
621                     return TableInsertAsLeft;
622                 }
623             }
624             else if (StartingVpn <= NodeToExamine->EndingVpn) {

625
626                 //
627                 // This is the node.
628                 //

629
630                 *NodeOrParent = NodeToExamine;
631                 return TableFoundNode;
632             }
633             else {

634
635                 Child = NodeToExamine->RightChild;

636
637                 if (Child != NULL) {
638                     NodeToExamine = Child;
639                 }
640                 else {

641
642                     //
643                     // Node is not in the tree.  Set the output
644                     // parameter to point to what would be its
```

```
645                          // parent and return which child it would be.
646                          //

647
648                          *NodeOrParent = NodeToExamine;
649                          return TableInsertAsRight;
650                      }
651                  }

652
653          } while (TRUE);
654  }

655

656
657  PMMADDRESS_NODE
658  MiCheckForConflictingNode (
659      IN ULONG_PTR StartVpn,
660      IN ULONG_PTR EndVpn,
661      IN PMM_AVL_TABLE Table
662      )

663
664  /*++

665
666  Routine Description:

667
668      The function determines if any addresses between a given starting
                and
669      ending address is contained within a virtual address descriptor.

670
671  Arguments:

672
673      StartVpn - Supplies the virtual address to locate a containing
674                          descriptor.

675
676      EndVpn - Supplies the virtual address to locate a containing
677                          descriptor.

678
679  Return Value:

680
681      Returns a pointer to the first conflicting virtual address
                descriptor
```

```
682         if one is found, otherwise a NULL value is returned.
683
684  ---*/
685
686  {
687      PMMADDRESS_NODE Node;
688
689      if (Table->NumberGenericTableElements == 0) {
690          return NULL;
691      }
692
693      Node = (PMMADDRESS_NODE) Table->BalancedRoot.RightChild;
694      ASSERT (Node != NULL);
695
696      do {
697
698          if (Node == NULL) {
699              return NULL;
700          }
701
702          if (StartVpn > Node->EndingVpn) {
703              Node = Node->RightChild;
704          }
705          else if (EndVpn < Node->StartingVpn) {
706              Node = Node->LeftChild;
707          }
708          else {
709
710              //
711              // The starting address is less than or equal to the end
                        VA
712              // and the ending address is greater than or equal to the
713              // start va.  Return this node.
714              //
715
716              return Node;
717          }
718
719      } while (TRUE);
```

```
720  }
721
722
723  PMMADDRESS_NODE
724  FASTCALL
725  MiGetFirstNode (
726      IN PMM_AVL_TABLE Table
727      )
728
729  /*++
730
731  Routine Description:
732
733      This function locates the virtual address descriptor which
                contains
734      the address range which logically is first within the address
                space.
735
736  Arguments:
737
738      None.
739
740  Return Value:
741
742      Returns a pointer to the virtual address descriptor containing
                the
743      first address range, NULL if none.
744
745  --*/
746
747  {
748      PMMADDRESS_NODE First;
749
750      if (Table->NumberGenericTableElements == 0) {
751          return NULL;
752      }
753
754      First = (PMMADDRESS_NODE) Table->BalancedRoot.RightChild;
755
```

```
756        ASSERT (First != NULL);

757

758        while (First->LeftChild != NULL) {

759            First = First->LeftChild;

760        }

761

762        return First;

763    }

764

765

766    VOID

767    MiPromoteNode (

768        IN PMMADDRESS_NODE C

769        )

770

771    /*++

772

773    Routine Description:

774

775        This routine performs the fundamental adjustment required for
               balancing

776        the binary tree during insert and delete operations.  Simply put,
                the

777        designated node is promoted in such a way that it rises one level
                in

778        the tree and its parent drops one level in the tree, becoming now
                the

779        child of the designated node.  Generally the path length to the
               subtree

780        "opposite" the original parent.  Balancing occurs as the caller
               chooses

781        which nodes to promote according to the balanced tree algorithms
               from

782        Knuth.

783

784        This is not the same as a splay operation, typically a splay "
               promotes"

785        a designated node twice.

786
```

```
787        Note that the pointer to the root node of the tree is assumed to
              be
788        contained in a MMADDRESS_NODE structure itself, to allow the
789        algorithms below to change the root of the tree without checking
790        for special cases.  Note also that this is an internal routine,
791        and the caller guarantees that it never requests to promote the
792        root itself.
793
794        This routine only updates the tree links; the caller must update
795        the balance factors as appropriate.
796
797  Arguments:
798
799        C - pointer to the child node to be promoted in the tree.
800
801  Return Value:
802
803        None.
804
805  --*/
806
807  {
808        PMMADDRESS_NODE P;
809        PMMADDRESS_NODE G;
810
811        //
812        // Capture the current parent and grandparent (may be the root).
813        //
814
815        P = SANITIZE_PARENT_NODE (C->u1.Parent);
816        G = SANITIZE_PARENT_NODE (P->u1.Parent);
817
818        //
819        // Break down the promotion into two cases based upon whether C
820        // is a left or right child.
821        //
822
823        if (P->LeftChild == C) {
824
```

```
825         //
826         // This promotion looks like this:
827         //
828         //            G                    G
829         //            |                    |
830         //            P                    C
831         //           / \     =>          / \
832         //          C   z               x   P
833         //         / \                     / \
834         //        x   y                   y   z
835         //

837         P->LeftChild = C->RightChild;

839         if (P->LeftChild != NULL) {

841             P->LeftChild->u1.Parent = MI_MAKE_PARENT (P, P->LeftChild
                    ->u1.Balance);
842         }

844         C->RightChild = P;

846         //
847         // Fall through to update parent and G <-> C relationship in
848         // common code.
849         //

851     }
852     else {

854         ASSERT(P->RightChild == C);

856         //
857         // This promotion looks like this:
858         //
859         //            G                    G
860         //            |                    |
861         //            P                    C
862         //           / \     =>          / \
```

```
863         //        x   C                P    z
864         //             / \              / \
865         //            y   z            x   y
866         //

868         P->RightChild = C->LeftChild;

870         if (P->RightChild != NULL) {
871             P->RightChild->u1.Parent = MI_MAKE_PARENT (P, P->
                    RightChild->u1.Balance);
872         }

874         C->LeftChild = P;
875     }

877     //
878     // Update parent of P, for either case above.
879     //

881     P->u1.Parent = MI_MAKE_PARENT (C, P->u1.Balance);

883     //
884     // Finally update G <-> C links for either case above.
885     //

887     if (G->LeftChild == P) {
888         G->LeftChild = C;
889     }
890     else {
891         ASSERT(G->RightChild == P);
892         G->RightChild = C;
893     }
894     C->u1.Parent = MI_MAKE_PARENT (G, C->u1.Balance);
895 }


898 ULONG
899 MiRebalanceNode (
900     IN PMMADDRESS_NODE S
```

```
901        )
902
903   /*++
904
905   Routine Description:
906
907       This routine performs a rebalance around the input node S, for
               which the
908       Balance factor has just effectively become +2 or −2.  When called
               , the
909       Balance factor still has a value of +1 or −1, but the respective
               longer
910       side has just become one longer as the result of an insert or
               delete
911       operation.
912
913       This routine effectively implements steps A7.iii (test for Case 1
                or
914       Case 2) and steps A8 and A9 of Knuth's balanced insertion
               algorithm,
915       plus it handles Case 3 identified in the delete section, which
               can
916       only happen on deletes.
917
918       The trick is, to convince yourself that while traveling from the
919       insertion point at the bottom of the tree up, that there are only
920       these two cases, and that when traveling up from the deletion
               point,
921       that there are just these three cases.  Knuth says it is obvious!
922
923   Arguments:
924
925       S − pointer to the node which has just become unbalanced.
926
927   Return Value:
928
929       TRUE if Case 3 was detected (causes delete algorithm to terminate
               ).
930
```

```
931  Environment:
932
933      Kernel mode.  The PFN lock is held for some of the tables.
934
935  --*/
936
937  {
938      PMMADDRESS_NODE R, P;
939      SCHAR a;
940
941      PRINT("rebalancing node %p bal=%x start=%x end=%x\n",
942                      S,
943                      S->u1.Balance,
944                      S->StartingVpn,
945                      S->EndingVpn);
946
947      //
948      // The parent node is never the argument node.
949      //
950
951      ASSERT (SANITIZE_PARENT_NODE(S->u1.Parent) != S);
952
953      //
954      // Capture which side is unbalanced.
955      //
956
957      a = (SCHAR) S->u1.Balance;
958
959      if (a == +1) {
960          R = S->RightChild;
961      }
962      else {
963          R = S->LeftChild;
964      }
965
966      //
967      // If the balance of R and S are the same (Case 1 in Knuth) then
            a single
968      // promotion of R will do the single rotation.  (Step A8, A10)
```

```
969        //
970        //  Here  is  a  diagram  of  the  Case  1  transformation ,  for  a ==  +1 ( a
               mirror
971        //  image  transformation  occurs  when  a == −1) ,  and  where  the
               subtree
972        //  heights  are  h  and  h+1  as  shown  (++  indicates  the  node  out  of
               balance ) :
973        //
974        //                           |                           |
975        //                          S++                          R
976        //                         / \                          / \
977        //                     ( h )   R+       ==>         S   ( h+1 )
978        //                          / \                    / \
979        //                      ( h ) ( h+1 )           ( h ) ( h )
980        //
981        //  Note  that  on  an  insert  we  can  hit  this  case  by  inserting  an
               item  in  the
982        //  right  subtree  of  R.   The  original  height  of  the  subtree  before
                the  insert
983        //  was  h+2,  and  it  is  still  h+2  after  the  rebalance ,  so  insert
               rebalancing
984        //  may  terminate .
985        //
986        //  On  a  delete  we  can  hit  this  case  by  deleting  a  node  from  the
               left  subtree
987        //  of  S.   The  height  of  the  subtree  before  the  delete  was  h+3,
               and  after  the
988        //  rebalance  it  is  h+2,  so  rebalancing  must  continue  up  the  tree .
989        //

991        if  ((SCHAR)  R−>u1 . Balance == a)  {

993            MiPromoteNode  (R) ;
994            R−>u1 . Balance  =  0;
995            S−>u1 . Balance  =  0;

997            return  FALSE;
998        }

999
```

```
//
// Otherwise, we have to promote the appropriate child of R twice (Case 2
// in Knuth).  (Step A9, A10)
//
// Here is a diagram of the Case 2 transformation, for a == +1 (a mirror
// image transformation occurs when a == -1), and where the subtree
// heights are h and h-1 as shown.  There are actually two minor subcases,
// differing only in the original balance of P (++ indicates the node out
// of balance).
//
//                      |                       |
//                     S++                      P
//                    / \                      / \
//                   /   \                    /   \
//                  /     \                  /     \
//               (h)      R-     ==>      S-        R
//                       / \            / \       / \
//                     P+ (h)        (h)(h-1)(h) (h)
//                    / \
//                 (h-1) (h)
//
//
//                      |                       |
//                     S++                      P
//                    / \                      / \
//                   /   \                    /   \
//                  /     \                  /     \
//               (h)      R-     ==>      S        R+
//                       / \            / \       / \
//                     P- (h)        (h) (h)(h-1)(h)
//                    / \
//                 (h) (h-1)
//
// Note that on an insert we can hit this case by inserting an
```
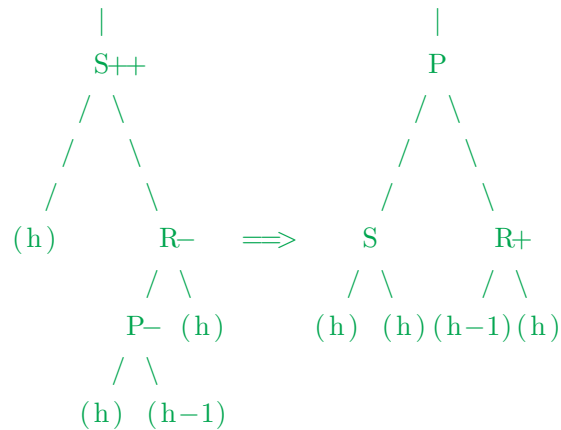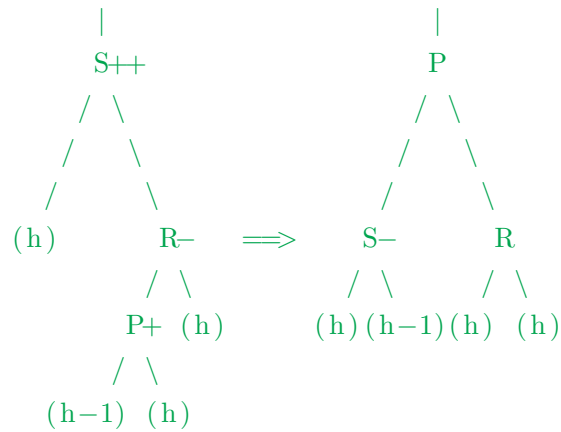
```
                       item in the
1034       // left subtree of R.  The original height of the subtree before
               the insert
1035       // was h+2, and it is still h+2 after the rebalance, so insert
               rebalancing
1036       // may terminate.
1037       //
1038       // On a delete we can hit this case by deleting a node from the
               left subtree
1039       // of S.  The height of the subtree before the delete was h+3,
               and after the
1040       // rebalance it is h+2, so rebalancing must continue up the tree.
1041       //
1042
1043       if ((SCHAR) R->u1.Balance == -a) {
1044
1045           //
1046           // Pick up the appropriate child P for the double rotation (
                   Link(-a,R)).
1047           //
1048
1049           if (a == 1) {
1050               P = R->LeftChild;
1051           }
1052           else {
1053               P = R->RightChild;
1054           }
1055
1056           //
1057           // Promote him twice to implement the double rotation.
1058           //
1059
1060           MiPromoteNode (P);
1061           MiPromoteNode (P);
1062
1063           //
1064           // Now adjust the balance factors.
1065           //
1066
```

```
1067            S->u1.Balance = 0;
1068            R->u1.Balance = 0;
1069            if ((SCHAR) P->u1.Balance == a) {
1070                PRINT("REBADJ_A:_Node_%p,_Bal_%x_->_%x\n", S, S->u1.
                        Balance, -a);
1071                COUNT_BALANCE_MAX ((SCHAR)-a);
1072                S->u1.Balance = (ULONG_PTR) -a;
1073            }
1074            else if ((SCHAR) P->u1.Balance == -a) {
1075                PRINT("REBADJ_B:_Node_%p,_Bal_%x_->_%x\n", R, R->u1.
                        Balance, a);
1076                COUNT_BALANCE_MAX ((SCHAR)a);
1077                R->u1.Balance = (ULONG_PTR) a;
1078            }
1079
1080            P->u1.Balance = 0;
1081            return FALSE;
1082        }
1083
1084        //
1085        // Otherwise this is Case 3 which can only happen on Delete (
                identical
1086        // to Case 1 except R->u1.Balance == 0).  We do a single rotation
                , adjust
1087        // the balance factors appropriately, and return TRUE.  Note that
                 the
1088        // balance of S stays the same.
1089        //
1090        // Here is a diagram of the Case 3 transformation, for a == +1 (a
                 mirror
1091        // image transformation occurs when a == -1), and where the
                subtree
1092        // heights are h and h+1 as shown (++ indicates the node out of
                balance):
1093        //
1094        //                     |                       |
1095        //                    S++                     R-
1096        //                   / \                     / \
1097        //                 (h)  R        ==>       S+  (h+1)
```

```
1098        //                          / \                / \
1099        //                       (h+1)(h+1)         (h)  (h+1)
1100        //
1101        // This case can not occur on an insert, because it is impossible
               for
1102        // a single insert to balance R, yet somehow grow the right
               subtree of
1103        // S at the same time.  As we move up the tree adjusting balance
               factors
1104        // after an insert, we terminate the algorithm if a node becomes
               balanced,
1105        // because that means the subtree length did not change!
1106        //
1107        // On a delete we can hit this case by deleting a node from the
               left
1108        // subtree of S.  The height of the subtree before the delete was
               h+3,
1109        // and after the rebalance it is still h+3, so rebalancing may
               terminate
1110        // in the delete path.
1111        //

1113        MiPromoteNode (R);
1114        PRINT("REBADJ_C:_Node_%p,_Bal_%x_->_%x\n", R, R->u1.Balance, -a);
1115        COUNT_BALANCE_MAX ((SCHAR)-a);
1116        R->u1.Balance = -a;
1117        return TRUE;
1118    }



1121    VOID
1122    FASTCALL
1123    MiRemoveNode (
1124        IN PMMADDRESS_NODE NodeToDelete,
1125        IN PMM_AVL_TABLE Table
1126        )

1128    /*++
1129
```

```
1130  Routine Description:
1131
1132      This routine deletes the specified node from the balanced tree,
              rebalancing
1133      as necessary.  If the NodeToDelete has at least one NULL child
              pointers,
1134      then it is chosen as the EasyDelete, otherwise a subtree
              predecessor or
1135      successor is found as the EasyDelete.  In either case the
              EasyDelete is
1136      deleted and the tree is rebalanced.  Finally if the NodeToDelete
              was
1137      different than the EasyDelete, then the EasyDelete is linked back
               into the
1138      tree in place of the NodeToDelete.
1139
1140  Arguments:
1141
1142      NodeToDelete - Pointer to the node which the caller wishes to
              delete.
1143
1144      Table - The generic table in which the delete is to occur.
1145
1146  Return Value:
1147
1148      None.
1149
1150  Environment:
1151
1152      Kernel mode.  The PFN lock is held for some of the tables.
1153
1154  --*/
1155
1156  {
1157      PMMADDRESS_NODE rebalance;
1158      Table->NumberGenericTableElements -= 1;
1159      rebalance = ___rb_erase_augmented(NodeToDelete, Table);
1160      if (rebalance!=&Table->BalancedRoot&&rebalance)
1161          _____rb_erase_color(rebalance, Table);
```

```
1162  }
1163
1164
1165  PMMADDRESS_NODE
1166  MiRealSuccessor (
1167      IN PMMADDRESS_NODE Links
1168      )
1169
1170  /*++
1171
1172  Routine Description:
1173
1174      This function takes as input a pointer to a balanced link
1175      in a tree and returns a pointer to the successor of the input
            node within
1176      the entire tree.  If there is not a successor, the return value
            is NULL.
1177
1178  Arguments:
1179
1180      Links − Supplies a pointer to a balanced link in a tree.
1181
1182  Return Value:
1183
1184      PMMADDRESS_NODE − returns a pointer to the successor in the
            entire tree
1185
1186  −−*/
1187
1188  {
1189      PMMADDRESS_NODE Ptr;
1190
1191      /*
1192          First check to see if there is a right subtree to the input
                link
1193          if there is then the real successor is the left most node in
1194          the right subtree.  That is find and return S in the
                following diagram
1195
```

```
1196                    Links
1197                        \
1198                          .
1199                          .
1200                          .
1201                        /
1202                    S
1203                        \
1204        */
1205
1206        if ((Ptr = Links->RightChild) != NULL) {
1207
1208            while (Ptr->LeftChild != NULL) {
1209                Ptr = Ptr->LeftChild;
1210            }
1211
1212            return Ptr;
1213        }
1214
1215        /*
1216            We do not have a right child so check to see if have a parent
                    and if
1217            so find the first ancestor that we are a left decendant of.
                    That
1218            is find and return S in the following diagram
1219
1220                        S
1221                      /
1222                    .
1223                    .
1224                    .
1225                        Links
1226
1227            Note that this code depends on how the BalancedRoot is
                    initialized,
1228            which is Parent points to self, and the RightChild points to
                    an
1229            actual node which is the root of the tree, and LeftChild does
                    not
```

```
            point  to  self .
        */

        Ptr = Links;
        while  ( MiIsRightChild ( Ptr ))  {
            Ptr = SANITIZE_PARENT_NODE  ( Ptr−>u1 . Parent );
        }

        if  ( MiIsLeftChild ( Ptr ))  {
            return  SANITIZE_PARENT_NODE  ( Ptr−>u1 . Parent );
        }

        //
        // Otherwise  we  are  do  not  have  a  real  successor  so  we  simply
            return  NULL.
        //
        // This  can  only  occur  when  we  get  back  to  the  root , and  we  can
            tell
        // that  since  the  Root  is  its  own  parent .
        //

        ASSERT  ( SANITIZE_PARENT_NODE( Ptr−>u1 . Parent )  ==  Ptr );

        return  NULL;
}


PMMADDRESS_NODE
MiRealPredecessor (
    IN PMMADDRESS_NODE Links
    )

/*++

Routine  Description :

    The  RealPredecessor  function  takes  as  input  a  pointer  to  a
        balanced  link
    in  a  tree  and  returns  a  pointer  to  the  predecessor  of  the  input
```

```
                node
1266    within the entire tree.  If there is not a predecessor, the
           return value
1267    is NULL.
1268
1269 Arguments:
1270
1271    Links - Supplies a pointer to a balanced link in a tree.
1272
1273 Return Value:
1274
1275    PMMADDRESS_NODE - returns a pointer to the predecessor in the
           entire tree
1276
1277 --*/
1278
1279 {
1280    PMMADDRESS_NODE Ptr;
1281    PMMADDRESS_NODE Parent;
1282    PMMADDRESS_NODE GrandParent;
1283
1284    /*
1285       First check to see if there is a left subtree to the input link
1286       if there is then the real predecessor is the right most node in
1287       the left subtree.  That is find and return P in the following
1288          diagram
1289                      Links
1290                       /
1291                      .
1292                       .
1293                        .
1294                       P
1295                      /
1296    */
1297
1298    if ((Ptr = Links->LeftChild) != NULL) {
1299
1300        while (Ptr->RightChild != NULL) {
```

```
1301              Ptr = Ptr->RightChild;
1302          }
1303
1304          return Ptr;
1305
1306      }
1307
1308      /*
1309        We do not have a left child so check to see if have a parent
                and if
1310        so find the first ancestor that we are a right decendant of.
                That
1311        is find and return P in the following diagram
1312
1313                              P
1314                               \
1315                                .
1316                                .
1317                                .
1318                            Links
1319
1320          Note that this code depends on how the BalancedRoot is
                    initialized ,
1321          which is Parent points to self , and the RightChild points to
                    an
1322          actual node which is the root of the tree .
1323      */
1324
1325      Ptr = Links ;
1326      while (MiIsLeftChild(Ptr)) {
1327          Ptr = SANITIZE_PARENT_NODE (Ptr->u1.Parent);
1328      }
1329
1330      if (MiIsRightChild(Ptr)) {
1331          Parent = SANITIZE_PARENT_NODE (Ptr->u1.Parent);
1332          GrandParent = SANITIZE_PARENT_NODE (Parent->u1.Parent);
1333          if (GrandParent != Parent) {
1334              return Parent;
1335          }
```

```
1336        }
1337
1338        //
1339        // Otherwise we are do not have a real predecessor so we simply
                   return
1340        // NULL.
1341        //
1342
1343        return NULL;
1344  }
1345
1346
1347  VOID
1348  MiInitializeVadTableAvl (
1349        IN PMM_AVL_TABLE Table
1350        )
1351
1352  /*++
1353
1354  Routine Description:
1355
1356        This routine initializes a table.
1357
1358  Arguments:
1359
1360        Table − Pointer to the generic table to be initialized.
1361
1362  Return Value:
1363
1364        None.
1365
1366  −−*/
1367
1368  {
1369
1370  #if DBG
1371        ULONG i;
1372
1373        for (i = 2; i < 33; i += 1) {
```

```
1374            ASSERT(MiWorstCaseFill[i] == (1 + MiWorstCaseFill[i − 1] +
                    MiWorstCaseFill[i − 2]));
1375        }
1376  #endif
1377
1378      //
1379      // Initialize each field in the argument Table.
1380      //
1381
1382      RtlZeroMemory (Table, sizeof(MM_AVL_TABLE));
1383
1384      Table−>BalancedRoot.u1.Parent = MI_MAKE_PARENT (&Table−>
              BalancedRoot, 0);
1385  }
1386
1387
1388  VOID
1389  FASTCALL
1390  MiInsertNode (
1391      IN PMMADDRESS_NODE NodeToInsert,
1392      IN PMM_AVL_TABLE Table
1393      )
1394
1395  /*++
1396
1397  Routine Description:
1398
1399      This function inserts a new element in a table.
1400
1401  Arguments:
1402
1403      NodeToInsert − The initialized address node to insert.
1404
1405      Table − Pointer to the table in which to insert the new node.
1406
1407  Return Value:
1408
1409      None.
1410
```

```
1411  Environment:
1412
1413      Kernel mode.  The PFN lock is held for some of the tables.
1414
1415  --*/
1416
1417  {
1418      PMMADDRESS_NODE NodeOrParent;
1419      PMMADDRESS_NODE parent,gparent,tmp;
1420      TABLE_SEARCH_RESULT SearchResult;
1421
1422      SearchResult = MiFindNodeOrParent (Table,
1423          NodeToInsert->StartingVpn,
1424          &NodeOrParent);
1425
1426      NodeToInsert->LeftChild = NULL;
1427      NodeToInsert->RightChild = NULL;
1428
1429      Table->NumberGenericTableElements += 1;
1430
1431      //
1432      // Insert the newer node in the tree.
1433      //
1434
1435      if (SearchResult == TableEmptyTree)
1436      {
1437          Table->BalancedRoot.RightChild = NodeToInsert;
1438          rb_set_parent(NodeToInsert,&Table->BalancedRoot);
1439      }
1440      else
1441      {
1442          if (SearchResult == TableInsertAsLeft)
1443          {
1444              NodeOrParent->LeftChild = NodeToInsert;
1445          }
1446          else
1447          {
1448              NodeOrParent->RightChild = NodeToInsert;
1449          }
```

```
1450            rb_set_parent(NodeToInsert,NodeOrParent);
1451        }
1452    rb_set_red(NodeToInsert);
1453    parent=rb_parent(NodeToInsert);
1454    while(1)   {
1455        if (parent==&Table->BalancedRoot) {
1456            Table->BalancedRoot.RightChild = NodeToInsert;
1457            rb_set_parent_color(NodeToInsert,&Table->BalancedRoot,
                    rb_black);
1458            break;
1459        } else if (rb_is_black(parent))
1460            break;
1461        gparent = rb_parent(parent);
1462        tmp = gparent->RightChild;
1463
1464        if (parent != tmp) {      /* parent == gparent->rb_left */
1465            if (tmp && rb_is_red(tmp)) {
1466                /*
1467                 * Case 1 - color flips
1468                 *
1469                 *      G              g
1470                 *     / \            / \
1471                 *    p   u  -->    P   U
1472                 *   /              /
1473                 *  n              n
1474                 *
1475                 * However, since g's parent might be red, and
1476                 * 4) does not allow this, we need to recurse
1477                 * at g.
1478                 */
1479                rb_set_parent_color(tmp, gparent, rb_black);
1480                rb_set_parent_color(parent, gparent, rb_black);
1481                NodeToInsert = gparent;
1482                parent = rb_parent(NodeToInsert);
1483                rb_set_parent_color(NodeToInsert, parent, rb_red);
1484                continue;
1485            }
1486
1487            tmp = parent->RightChild;
```

```
1488            if (NodeToInsert == tmp) {
1489                /*
1490                 * Case 2 - left rotate at parent
1491                 *
1492                 *      G                G
1493                 *     / \              / \
1494                 *    p   U   -->     n   U
1495                 *     \              /
1496                 *      n            p
1497                 *
1498                 * This still leaves us in violation of 4), the
1499                 * continuation into Case 3 will fix that.
1500                 */
1501                parent->RightChild = tmp = NodeToInsert->LeftChild;
1502                NodeToInsert->LeftChild = parent;
1503                if (tmp)
1504                    rb_set_parent_color(tmp, parent,rb_black);
1505                rb_set_parent_color(parent, NodeToInsert, rb_red);
1506                parent = NodeToInsert;
1507                tmp = NodeToInsert->RightChild;
1508            }
1509
1510            /*
1511             * Case 3 - right rotate at gparent
1512             *
1513             *        G              P
1514             *       / \            / \
1515             *      p   U   -->  n   g
1516             *     /                  \
1517             *    n                    U
1518             */
1519            gparent->LeftChild= tmp;   /* == parent->rb_right */
1520            parent->RightChild= gparent;
1521            if (tmp)
1522                rb_set_parent_color(tmp, gparent, rb_black);
1523            __rb_rotate_set_parents(gparent, parent, Table, rb_red);
1524            break;
1525        } else {
1526            tmp = gparent->LeftChild;
```

```
1527                    if  (tmp && rb_is_red(tmp))  {
1528                        /* Case  1 – color  flips  */
1529                        rb_set_parent_color(tmp, gparent, rb_black);
1530                        rb_set_parent_color(parent, gparent, rb_black);
1531                        NodeToInsert = gparent;
1532                        parent = rb_parent(NodeToInsert);
1533                        rb_set_parent_color(NodeToInsert, parent, rb_red);
1534                        continue;
1535                    }
1536
1537                    tmp = parent->LeftChild;
1538                    if  (NodeToInsert == tmp)  {
1539                        /* Case  2 – right  rotate  at  parent  */
1540                        parent->LeftChild = tmp = NodeToInsert->RightChild;
1541                        NodeToInsert->RightChild = parent;
1542                        if  (tmp)
1543                            rb_set_parent_color(tmp, parent,
1544                            rb_black);
1545                        rb_set_parent_color(parent, NodeToInsert, rb_red);
1546                        parent = NodeToInsert;
1547                        tmp = NodeToInsert->LeftChild;
1548                    }
1549
1550                    /* Case  3 – left  rotate  at  gparent  */
1551                    gparent->RightChild = tmp;   /* == parent->rb_left  */
1552                    parent->LeftChild = gparent;
1553                    if  (tmp)
1554                        rb_set_parent_color(tmp, gparent, rb_black);
1555                    __rb_rotate_set_parents(gparent, parent, Table, rb_red);
1556                    break;
1557                }
1558            }
1559        return;
1560    }
1561
1562
1563    PVOID
1564    MiEnumerateGenericTableWithoutSplayingAvl (
1565        IN PMM_AVL_TABLE Table,
```

```
1566      IN PVOID *RestartKey
1567      )
1568
1569  /*++
1570
1571  Routine Description:
1572
1573      The function EnumerateGenericTableWithoutSplayingAvl will return
              to the
1574      caller one-by-one the elements of of a table.  The return value
              is a
1575      pointer to the user defined structure associated with the element
              .
1576      The input parameter RestartKey indicates if the enumeration
              should
1577      start from the beginning or should return the next element.  If
              the
1578      are no more new elements to return the return value is NULL.  As
              an
1579      example of its use, to enumerate all of the elements in a table
              the
1580      user would write:
1581
1582          *RestartKey = NULL;
1583
1584          for (ptr = EnumerateGenericTableWithoutSplayingAvl(Table, &
              RestartKey);
1585              ptr != NULL;
1586              ptr = EnumerateGenericTableWithoutSplayingAvl(Table, &
                  RestartKey)) {
1587                  :
1588          }
1589
1590  Arguments:
1591
1592      Table - Pointer to the generic table to enumerate.
1593
1594      RestartKey - Pointer that indicates if we should restart or
              return the next
```

```
1595                        element.    If  the  contents  of  RestartKey  is  NULL,  the
                                    search
1596                        will  be  started  from  the  beginning.
1597
1598  Return Value:
1599
1600      PVOID − Pointer  to  the  user  data.
1601
1602  −−*/
1603
1604  {
1605      PMMADDRESS_NODE NodeToReturn;
1606
1607      if (Table−>NumberGenericTableElements == 0) {
1608
1609          //
1610          // Nothing  to  do  if  the  table  is  empty.
1611          //
1612
1613          return  NULL;
1614
1615      }
1616
1617      //
1618      // If  the  restart  flag  is  true  then  go  to  the  least  element
1619      // in  the  tree.
1620      //
1621
1622      if (*RestartKey == NULL) {
1623
1624          //
1625          // Loop  until  we  find  the  leftmost  child  of  the  root.
1626          //
1627
1628          for  (NodeToReturn = Table−>BalancedRoot.RightChild;
1629                NodeToReturn−>LeftChild;
1630                NodeToReturn = NodeToReturn−>LeftChild) {
1631
1632              NOTHING;
```

```
1633            }

1634

1635            *RestartKey = NodeToReturn;

1636

1637        }
1638        else {

1639

1640            //
1641            // The caller has passed in the previous entry found
1642            // in the table to enable us to continue the search.  We call
1643            // RealSuccessor to step to the next element in the tree.
1644            //

1645

1646            NodeToReturn = MiRealSuccessor (*RestartKey);

1647

1648            if (NodeToReturn) {
1649                *RestartKey = NodeToReturn;
1650            }
1651        }

1652

1653        //
1654        // Return the found element.
1655        //

1656

1657        return NodeToReturn;
1658  }

1659

1660

1661  PMMADDRESS_NODE
1662  FASTCALL
1663  MiGetNextNode (
1664        IN PMMADDRESS_NODE Node
1665        )

1666

1667  /*++

1668

1669  Routine Description:

1670

1671        This function locates the virtual address descriptor which
```

```
                     contains
1672       the address range which logically follows the specified address
              range.

1673

1674   Arguments:

1675

1676       Node - Supplies a pointer to a virtual address descriptor.

1677

1678   Return Value:

1679

1680       Returns a pointer to the virtual address descriptor containing
              the
1681       next address range, NULL if none.

1682

1683   --*/

1684

1685   {

1686       PMMADDRESS_NODE Next;

1687       PMMADDRESS_NODE Parent;

1688       PMMADDRESS_NODE Left;

1689

1690       Next = Node;

1691

1692       if (Next->RightChild == NULL) {

1693

1694           do {

1695

1696               Parent = SANITIZE_PARENT_NODE (Next->u1.Parent);

1697

1698               ASSERT (Parent != NULL);

1699

1700               if (Parent == Next) {

1701                   return NULL;

1702               }

1703

1704               //
1705               // Locate the first ancestor of this node of which this
1706               // node is the left child of and return that node as the
1707               // next element.
```

```
1708                    //
1709
1710                    if (Parent->LeftChild == Next) {
1711                        return Parent;
1712                    }
1713
1714                    Next = Parent;
1715
1716            } while (TRUE);
1717        }
1718
1719        //
1720        // A right child exists, locate the left most child of that right
                  child.
1721        //
1722
1723        Next = Next->RightChild;
1724
1725        do {
1726
1727            Left = Next->LeftChild;
1728
1729            if (Left == NULL) {
1730                break;
1731            }
1732
1733            Next = Left;
1734
1735        } while (TRUE);
1736
1737        return Next;
1738
1739 }
1740
1741 PMMADDRESS_NODE
1742 FASTCALL
1743 MiGetPreviousNode (
1744     IN PMMADDRESS_NODE Node
1745     )
```

```
/*++

Routine Description:

    This function locates the virtual address descriptor which
        contains
    the address range which logically precedes the specified virtual
    address descriptor.

Arguments:

    Node - Supplies a pointer to a virtual address descriptor.

Return Value:

    Returns a pointer to the virtual address descriptor containing
        the
    next address range, NULL if none.

--*/

{
    PMMADDRESS_NODE Previous;
    PMMADDRESS_NODE Parent;

    Previous = Node;

    if (Previous->LeftChild == NULL) {

        ASSERT (Previous->u1.Parent != NULL);

        Parent = SANITIZE_PARENT_NODE (Previous->u1.Parent);

        while (Parent != Previous) {

            //
            // Locate the first ancestor of this node of which this
            // node is the right child of and return that node as the
```

```
1783              // Previous element.
1784              //
1785
1786              if (Parent->RightChild == Previous) {
1787
1788                  if (Parent == SANITIZE_PARENT_NODE (Parent->u1.Parent
                         )) {
1789                      return NULL;
1790                  }
1791
1792                  return Parent;
1793              }
1794
1795              Previous = Parent;
1796              Parent = SANITIZE_PARENT_NODE (Previous->u1.Parent);
1797          }
1798          return NULL;
1799      }
1800
1801      //
1802      // A left child exists, locate the right most child of that left
             child.
1803      //
1804
1805      Previous = Previous->LeftChild;
1806
1807      while (Previous->RightChild != NULL) {
1808          Previous = Previous->RightChild;
1809      }
1810
1811      return Previous;
1812  }
1813
1814
1815  PMMADDRESS_NODE
1816  FASTCALL
1817  MiLocateAddressInTree (
1818      IN ULONG_PTR Vpn,
1819      IN PMM_AVL_TABLE Table
```

```
1820        )

1821

1822   /*++

1823

1824   Routine Description:

1825

1826       The function locates the virtual address descriptor which
                  describes

1827       a given address.

1828

1829   Arguments:

1830

1831       Vpn - Supplies the virtual page number to locate a descriptor for
                  .

1832

1833   Return Value:

1834

1835       Returns a pointer to the virtual address descriptor which
                  contains

1836       the supplied virtual address or NULL if none was located.

1837

1838   --*/

1839

1840   {

1841       PVOID NodeOrParent;

1842       TABLE_SEARCH_RESULT SearchResult;

1843

1844       //

1845       // Lookup the element and save the result.

1846       //

1847

1848       SearchResult = MiFindNodeOrParent (Table,

1849                                          Vpn,

1850                                          (PMMADDRESS_NODE *) &
                                                NodeOrParent);

1851

1852       if (SearchResult == TableFoundNode) {

1853

1854           //
```

```
1855          // Return the VAD.
1856          //
1857
1858          return (PMMADDRESS_NODE) NodeOrParent;
1859      }
1860
1861      return NULL;
1862  }
1863
1864
1865  NTSTATUS
1866  MiFindEmptyAddressRangeInTree (
1867      IN SIZE_T SizeOfRange,
1868      IN ULONG_PTR Alignment,
1869      IN PMM_AVL_TABLE Table,
1870      OUT PMMADDRESS_NODE *PreviousVad,
1871      OUT PVOID *Base
1872      )
1873
1874  /*++
1875
1876  Routine Description:
1877
1878      The function examines the virtual address descriptors to locate
1879      an unused range of the specified size and returns the starting
1880      address of the range.
1881
1882  Arguments:
1883
1884      SizeOfRange - Supplies the size in bytes of the range to locate.
1885
1886      Alignment - Supplies the alignment for the address.  Must be
1887                  a power of 2 and greater than the page_size.
1888
1889      Table - Supplies the root of the tree to search through.
1890
1891      PreviousVad - Supplies the Vad which is before this the found
1892                  address range.
1893
```

```
1894        Base - Receives the starting address of a suitable range on
                success.

1895

1896   Return Value:

1897

1898       NTSTATUS.

1899

1900   --*/

1901

1902   {

1903       PMMADDRESS_NODE Node;

1904       PMMADDRESS_NODE NextNode;

1905       ULONG_PTR AlignmentVpn;

1906       ULONG_PTR SizeOfRangeVpn;

1907

1908       AlignmentVpn = Alignment >> PAGE_SHIFT;

1909

1910       //

1911       // Locate the node with the lowest starting address.

1912       //

1913

1914       ASSERT (SizeOfRange != 0);

1915       SizeOfRangeVpn = (SizeOfRange + (PAGE_SIZE - 1)) >> PAGE_SHIFT;

1916       ASSERT (SizeOfRangeVpn != 0);

1917

1918       if (Table->NumberGenericTableElements == 0) {

1919           *Base = MM_LOWEST_USER_ADDRESS;

1920           return STATUS_SUCCESS;

1921       }

1922

1923       Node = Table->BalancedRoot.RightChild;

1924

1925       while (Node->LeftChild != NULL) {

1926           Node = Node->LeftChild;

1927       }

1928

1929       //

1930       // Check to see if a range exists between the lowest address VAD

1931       // and lowest user address.
```

```
1932        //
1933
1934        if (Node->StartingVpn > MI_VA_TO_VPN (MM_LOWEST_USER_ADDRESS)) {
1935
1936            if (SizeOfRangeVpn <
1937                (Node->StartingVpn - MI_VA_TO_VPN (MM_LOWEST_USER_ADDRESS
                        ))) {
1938
1939                *PreviousVad = NULL;
1940                *Base = MM_LOWEST_USER_ADDRESS;
1941                return STATUS_SUCCESS;
1942            }
1943        }
1944
1945        do {
1946
1947            NextNode = MiGetNextNode (Node);
1948
1949            if (NextNode != NULL) {
1950
1951                if (SizeOfRangeVpn <=
1952                    ((ULONG_PTR)NextNode->StartingVpn -
1953                                    MI_ROUND_TO_SIZE(1 + Node->EndingVpn,
1954                                            AlignmentVpn))) {
1955
1956                    //
1957                    // Check to ensure that the ending address aligned
                            upwards
1958                    // is not greater than the starting address.
1959                    //
1960
1961                    if ((ULONG_PTR)NextNode->StartingVpn >
1962                            MI_ROUND_TO_SIZE(1 + Node->EndingVpn,
1963                                        AlignmentVpn)) {
1964
1965                        *PreviousVad = Node;
1966                        *Base = (PVOID) MI_ROUND_TO_SIZE(
1967                                    (ULONG_PTR)MI_VPN_TO_VA_ENDING(Node->
                                        EndingVpn),
```

```
1968                                    Alignment);
1969                       return STATUS_SUCCESS;
1970                   }
1971               }
1972
1973          } else {
1974
1975              //
1976              // No more descriptors, check to see if this fits into
                     the remainder
1977              // of the address space.
1978              //
1979
1980              if ((((ULONG_PTR)Node->EndingVpn + MI_VA_TO_VPN(X64K)) <
1981                      MI_VA_TO_VPN (MM_HIGHEST_VAD_ADDRESS))
1982                          &&
1983                  (SizeOfRange <=
1984                      ((ULONG_PTR)MM_HIGHEST_VAD_ADDRESS -
1985                          (ULONG_PTR)MI_ROUND_TO_SIZE(
1986                          (ULONG_PTR)MI_VPN_TO_VA(Node->EndingVpn),
                                Alignment)))) {
1987
1988                  *PreviousVad = Node;
1989                  *Base = (PVOID) MI_ROUND_TO_SIZE(
1990                              (ULONG_PTR)MI_VPN_TO_VA_ENDING(Node->
                                    EndingVpn),
1991                                   Alignment);
1992                  return STATUS_SUCCESS;
1993              }
1994              return STATUS_NO_MEMORY;
1995          }
1996          Node = NextNode;
1997
1998      } while (TRUE);
1999  }
2000
2001  NTSTATUS
2002  MiFindEmptyAddressRangeDownTree (
2003      IN SIZE_T SizeOfRange,
```

```
2004        IN PVOID HighestAddressToEndAt,
2005        IN ULONG_PTR Alignment,
2006        IN PMM_AVL_TABLE Table,
2007        OUT PVOID *Base
2008        )

2009
2010   /*++

2011
2012   Routine Description:

2013
2014        The function examines the virtual address descriptors to locate
2015        an unused range of the specified size and returns the starting
2016        address of the range.  The function examines from the high
2017        addresses down and ensures that starting address is less than
2018        the specified address.

2019
2020        Note this cannot be used for the based section tree because only
2021        the nodes in that tree are stored as VAs instead of VPNs.

2022
2023   Arguments:

2024
2025        SizeOfRange - Supplies the size in bytes of the range to locate.

2026
2027        HighestAddressToEndAt - Supplies the virtual address that limits
2028                                the value of the ending address.  The
                                    ending
2029                                address of the located range must be less
2030                                than this address.

2031
2032        Alignment - Supplies the alignment for the address.  Must be
2033                    a power of 2 and greater than the page_size.

2034
2035        Table - Supplies the root of the tree to search through.

2036
2037        Base - Receives the starting address of a suitable range on
                   success.

2038
2039   Return Value:

2040
```

```
2041        NTSTATUS.

2042

2043  --*/

2044

2045  {

2046        PMMADDRESS_NODE Node;

2047        PMMADDRESS_NODE PreviousNode;

2048        ULONG_PTR AlignedEndingVa;

2049        PVOID OptimalStart;

2050        ULONG_PTR OptimalStartVpn;

2051        ULONG_PTR HighestVpn;

2052        ULONG_PTR AlignmentVpn;

2053

2054        //

2055        // Note this cannot be used for the based section tree because
              only

2056        // the nodes in that tree are stored as VAs instead of VPNs.

2057        //

2058

2059        ASSERT (Table != &MmSectionBasedRoot);

2060

2061        SizeOfRange = MI_ROUND_TO_SIZE (SizeOfRange, PAGE_SIZE);

2062

2063        if (((ULONG_PTR)HighestAddressToEndAt + 1) < SizeOfRange) {

2064            return STATUS_NO_MEMORY;

2065        }

2066

2067        ASSERT (HighestAddressToEndAt != NULL);

2068        ASSERT (HighestAddressToEndAt <= (PVOID)((ULONG_PTR)
              MM_HIGHEST_VAD_ADDRESS + 1));

2069

2070        HighestVpn = MI_VA_TO_VPN (HighestAddressToEndAt);

2071

2072        //

2073        // Locate the Node with the highest starting address.

2074        //

2075

2076        OptimalStart = (PVOID)(MI_ALIGN_TO_SIZE(

2077                              (((ULONG_PTR)HighestAddressToEndAt + 1) -
```

```
                                  SizeOfRange),
2078                                 Alignment));
2079
2080        if (Table->NumberGenericTableElements == 0) {
2081
2082            //
2083            // The tree is empty, any range is okay.
2084            //
2085
2086            *Base = OptimalStart;
2087            return STATUS_SUCCESS;
2088        }
2089
2090        Node = (PMMADDRESS_NODE) Table->BalancedRoot.RightChild;
2091
2092        //
2093        // See if an empty slot exists to hold this range, locate the
                largest
2094        // element in the tree.
2095        //
2096
2097        while (Node->RightChild != NULL) {
2098            Node = Node->RightChild;
2099        }
2100
2101        //
2102        // Check to see if a range exists between the highest address VAD
2103        // and the highest address to end at.
2104        //
2105
2106        AlignedEndingVa = (ULONG_PTR)MI_ROUND_TO_SIZE ((ULONG_PTR)
                MI_VPN_TO_VA_ENDING (Node->EndingVpn),
2107                                                 Alignment);
2108
2109        if (AlignedEndingVa < (ULONG_PTR)HighestAddressToEndAt) {
2110
2111            if ( SizeOfRange < ((ULONG_PTR)HighestAddressToEndAt -
                    AlignedEndingVa)) {
2112
```

```
2113              *Base = MI_ALIGN_TO_SIZE(
2114                                    ((ULONG_PTR)HighestAddressToEndAt −
                                            SizeOfRange),
2115                                    Alignment);
2116            return STATUS_SUCCESS;
2117        }
2118      }
2119
2120      //
2121      // Walk the tree backwards looking for a fit.
2122      //
2123
2124      OptimalStartVpn = MI_VA_TO_VPN (OptimalStart);
2125      AlignmentVpn = MI_VA_TO_VPN (Alignment);
2126
2127      do {
2128
2129          PreviousNode = MiGetPreviousNode (Node);
2130
2131          if (PreviousNode != NULL) {
2132
2133              //
2134              // Is the ending Va below the top of the address to end
                       at.
2135              //
2136
2137              if (PreviousNode−>EndingVpn < OptimalStartVpn) {
2138                  if ((SizeOfRange >> PAGE_SHIFT) <=
2139                      ((ULONG_PTR)Node−>StartingVpn −
2140                      (ULONG_PTR)MI_ROUND_TO_SIZE(1 + PreviousNode−>
                            EndingVpn,
2141                                              AlignmentVpn))) {
2142
2143                      //
2144                      // See if the optimal start will fit between
                            these
2145                      // two VADs.
2146                      //
2147
```

```
2148                        if ((OptimalStartVpn > PreviousNode->EndingVpn)
                              &&
2149                             (HighestVpn < Node->StartingVpn)) {
2150                             *Base = OptimalStart;
2151                             return STATUS_SUCCESS;
2152                        }
2153
2154                        //
2155                        // Check to ensure that the ending address
                              aligned upwards
2156                        // is not greater than the starting address.
2157                        //
2158
2159                        if ((ULONG_PTR)Node->StartingVpn >
2160                                (ULONG_PTR)MI_ROUND_TO_SIZE(1 +
                                        PreviousNode->EndingVpn,
2161                                                    AlignmentVpn)) {
2162
2163                            *Base = MI_ALIGN_TO_SIZE(
2164                                            (ULONG_PTR)MI_VPN_TO_VA (
                                                Node->StartingVpn) -
                                                SizeOfRange,
2165                                            Alignment);
2166                            return STATUS_SUCCESS;
2167                        }
2168                    }
2169                }
2170        } else {
2171
2172            //
2173            // No more descriptors, check to see if this fits into
                  the remainder
2174            // of the address space.
2175            //
2176
2177            if (Node->StartingVpn > MI_VA_TO_VPN (
                  MM_LOWEST_USER_ADDRESS)) {
2178                if ((SizeOfRange >> PAGE_SHIFT) <=
2179                        ((ULONG_PTR)Node->StartingVpn - MI_VA_TO_VPN (
```

```
                                    MM_LOWEST_USER_ADDRESS))) {

2180
2181                      //
2182                      // See if the optimal start will fit between
                             these
2183                      // two VADs.
2184                      //
2185
2186                      if (HighestVpn < Node->StartingVpn) {
2187                          *Base = OptimalStart;
2188                          return STATUS_SUCCESS;
2189                      }
2190
2191                      *Base = MI_ALIGN_TO_SIZE(
2192                                  (ULONG_PTR)MI_VPN_TO_VA (Node->
                                         StartingVpn) - SizeOfRange,
2193                                  Alignment);
2194                      return STATUS_SUCCESS;
2195                  }
2196              }
2197              return STATUS_NO_MEMORY;
2198          }
2199          Node = PreviousNode;
2200
2201      } while (TRUE);
2202  }
2203
2204
2205  NTSTATUS
2206  MiFindEmptyAddressRangeDownBasedTree (
2207      IN SIZE_T SizeOfRange,
2208      IN PVOID HighestAddressToEndAt,
2209      IN ULONG_PTR Alignment,
2210      IN PMM_AVL_TABLE Table,
2211      OUT PVOID *Base
2212      )
2213
2214  /*++
2215
```

```
Routine Description:

    The function examines the virtual address descriptors to locate
    an unused range of the specified size and returns the starting
    address of the range.  The function examines from the high
    addresses down and ensures that starting address is less than
    the specified address.

    Note this is only used for the based section tree because only
    the nodes in that tree are stored as VAs instead of VPNs.

Arguments:

    SizeOfRange - Supplies the size in bytes of the range to locate.

    HighestAddressToEndAt - Supplies the virtual address that limits
                            the value of the ending address.  The
                                 ending
                            address of the located range must be less
                            than this address.

    Alignment - Supplies the alignment for the address.  Must be
                a power of 2 and greater than the page_size.

    Table - Supplies the root of the tree to search through.

    Base - Receives the starting address of a suitable range on
           success.

Return Value:

    NTSTATUS.

--*/

{
    PMMADDRESS_NODE Node;
    PMMADDRESS_NODE PreviousNode;
    ULONG_PTR AlignedEndingVa;
```

```
2253        ULONG_PTR OptimalStart;

2254

2255        //
2256        // Note this is only used for the based section tree because only
2257        // the nodes in that tree are stored as VAs instead of VPNs.
2258        //

2259

2260        ASSERT (Table == &MmSectionBasedRoot);

2261

2262        SizeOfRange = MI_ROUND_TO_SIZE (SizeOfRange, PAGE_SIZE);

2263

2264        if (((ULONG_PTR)HighestAddressToEndAt + 1) < SizeOfRange) {
2265            return STATUS_NO_MEMORY;
2266        }

2267

2268        ASSERT (HighestAddressToEndAt != NULL);
2269        ASSERT (HighestAddressToEndAt <= (PVOID)((ULONG_PTR)
                MM_HIGHEST_VAD_ADDRESS + 1));

2270

2271        //
2272        // Locate the node with the highest starting address.
2273        //

2274

2275        OptimalStart = (ULONG_PTR) MI_ALIGN_TO_SIZE (
2276                                  (((ULONG_PTR)HighestAddressToEndAt + 1) -
                                    SizeOfRange),
2277                                  Alignment);

2278

2279        if (Table->NumberGenericTableElements == 0) {

2280

2281            //
2282            // The tree is empty, any range is okay.
2283            //

2284

2285            *Base = (PVOID) OptimalStart;
2286            return STATUS_SUCCESS;
2287        }

2288

2289        Node = (PMMADDRESS_NODE) Table->BalancedRoot.RightChild;
```

```
2290
2291        //
2292        // See if an empty slot exists to hold this range, locate the
                  largest
2293        // element in the tree.
2294        //
2295
2296        while (Node->RightChild != NULL) {
2297            Node = Node->RightChild;
2298        }
2299
2300        //
2301        // Check to see if a range exists between the highest address VAD
2302        // and the highest address to end at.
2303        //
2304
2305        AlignedEndingVa = MI_ROUND_TO_SIZE (Node->EndingVpn, Alignment);
2306
2307        PRINT("search down0: %p %p %p\n", AlignedEndingVa,
                  HighestAddressToEndAt, SizeOfRange);
2308
2309        if ((AlignedEndingVa < (ULONG_PTR)HighestAddressToEndAt) &&
2310            (SizeOfRange < ((ULONG_PTR)HighestAddressToEndAt -
                    AlignedEndingVa))) {
2311
2312            *Base = MI_ALIGN_TO_SIZE(
2313                                    ((ULONG_PTR)HighestAddressToEndAt -
                                        SizeOfRange),
2314                                    Alignment);
2315            return STATUS_SUCCESS;
2316        }
2317
2318        //
2319        // Walk the tree backwards looking for a fit.
2320        //
2321
2322        do {
2323
2324            PreviousNode = MiGetPreviousNode (Node);
```

```
2325
2326            PRINT("search␣down1:␣%p␣%p␣%p␣%p\n", PreviousNode, Node,
                    OptimalStart, Alignment);
2327
2328            if (PreviousNode == NULL) {
2329                break;
2330            }
2331
2332            //
2333            // Is the ending Va below the top of the address to end at.
2334            //
2335
2336            if (PreviousNode->EndingVpn < OptimalStart) {
2337
2338                if (SizeOfRange <= (Node->StartingVpn -
2339                        MI_ROUND_TO_SIZE(1 + PreviousNode->EndingVpn,
                            Alignment))) {
2340
2341                    //
2342                    // See if the optimal start will fit between these
                        two VADs.
2343                    //
2344
2345                    if ((OptimalStart > PreviousNode->EndingVpn) &&
2346                        ((ULONG_PTR) HighestAddressToEndAt < Node->
                            StartingVpn)) {
2347                        *Base = (PVOID) OptimalStart;
2348                        return STATUS_SUCCESS;
2349                    }
2350
2351                    //
2352                    // Check to ensure that the ending address aligned
                        upwards
2353                    // is not greater than the starting address.
2354                    //
2355
2356                    if (Node->StartingVpn >
2357                        MI_ROUND_TO_SIZE(1 + PreviousNode->EndingVpn,
                            Alignment)) {
```

```
2358
2359                        *Base = MI_ALIGN_TO_SIZE (Node−>StartingVpn −
                               SizeOfRange,
2360                                                  Alignment);
2361
2362                        return STATUS_SUCCESS;
2363                    }
2364                }
2365            }
2366
2367            Node = PreviousNode;
2368
2369        } while (TRUE);
2370
2371
2372        //
2373        // No more descriptors, check to see if this fits into the
               remainder
2374        // of the address space.
2375        //
2376
2377        if (Node−>StartingVpn > (ULONG_PTR) MM_LOWEST_USER_ADDRESS) {
2378
2379            if (SizeOfRange <= (Node−>StartingVpn − (ULONG_PTR)
                   MM_LOWEST_USER_ADDRESS)) {
2380
2381                //
2382                // See if the optimal start will fit between these two
                       VADs.
2383                //
2384
2385                if ((ULONG_PTR) HighestAddressToEndAt < Node−>StartingVpn
                       ) {
2386                    *Base = (PVOID) OptimalStart;
2387                    return STATUS_SUCCESS;
2388                }
2389
2390                *Base = MI_ALIGN_TO_SIZE (Node−>StartingVpn − SizeOfRange
                       ,
```

```
2391                                                  Alignment);
2392
2393                 return STATUS_SUCCESS;
2394             }
2395         }
2396         return STATUS_NO_MEMORY;
2397 }
2398
2399 #if !defined (_USERMODE)
2400
2401 PMMVAD
2402 FASTCALL
2403 MiLocateAddress (
2404     IN PVOID VirtualAddress
2405     )
2406
2407 /*++
2408
2409 Routine Description:
2410
2411     The function locates the virtual address descriptor which
2412             describes
2412     a given address.
2413
2414 Arguments:
2415
2416     VirtualAddress - Supplies the virtual address to locate a
2417             descriptor for.
2417
2418     Table - Supplies the table describing the tree.
2419
2420 Return Value:
2421
2422     Returns a pointer to the virtual address descriptor which
2423             contains
2423     the supplied virtual address or NULL if none was located.
2424
2425 --*/
2426
```

```
2427  {
2428      PMMVAD FoundVad;
2429      ULONG_PTR Vpn;
2430      PMM_AVL_TABLE Table;
2431      TABLE_SEARCH_RESULT SearchResult;
2432
2433      Table = &PsGetCurrentProcess ()->VadRoot;
2434
2435      //
2436      // Note the NodeHint *MUST* be captured locally - see the
             synchronization
2437      // comment below for details.
2438      //
2439
2440      FoundVad = (PMMVAD) Table->NodeHint;
2441
2442      if (FoundVad == NULL) {
2443          return NULL;
2444      }
2445
2446      Vpn = MI_VA_TO_VPN (VirtualAddress);
2447
2448      if ((Vpn >= FoundVad->StartingVpn) && (Vpn <= FoundVad->EndingVpn
             )) {
2449          return FoundVad;
2450      }
2451
2452      //
2453      // Lookup the element and save the result.
2454      //
2455
2456      SearchResult = MiFindNodeOrParent (Table,
2457                                         Vpn,
2458                                         (PMMADDRESS_NODE *) &FoundVad)
                                             ;
2459
2460      if (SearchResult != TableFoundNode) {
2461          return NULL;
2462      }
```

```
2463
2464        ASSERT (FoundVad != NULL);
2465
2466        ASSERT ((Vpn >= FoundVad->StartingVpn) && (Vpn <= FoundVad->
                EndingVpn));
2467
2468        //
2469        // Note the NodeHint field update is not synchronized in all
                cases, ie:
2470        // some callers hold the address space mutex and others hold the
                working
2471        // set pushlock.  It is ok that the update is not synchronized -
                as long
2472        // as care is taken above that it is read into a local variable
                and then
2473        // referenced.  Because no VAD can be removed from the tree
                without holding
2474        // both the address space & working set.
2475        //
2476
2477        Table->NodeHint = (PVOID) FoundVad;
2478
2479        //
2480        // Return the VAD.
2481        //
2482
2483        return FoundVad;
2484 }
2485 #endif
2486
2487 #if DBG
2488 VOID
2489 MiNodeTreeWalk (
2490        IN PMM_AVL_TABLE Table
2491        )
2492 {
2493        PVOID RestartKey;
2494        PMMADDRESS_NODE NewNode;
2495        PMMADDRESS_NODE PrevNode;
```

```
2496        PMMADDRESS_NODE NextNode;

2497

2498      RestartKey = NULL;

2499

2500      do {

2501

2502          NewNode = MiEnumerateGenericTableWithoutSplayingAvl (Table,
2503                                                                 &
                                                                     RestartKey
                                                                     );

2504

2505          if (NewNode == NULL) {
2506              break;
2507          }

2508

2509          PrevNode = MiGetPreviousNode (NewNode);
2510          NextNode = MiGetNextNode (NewNode);

2511

2512          PRINT ("Node %p %x %x\n",
2513                          NewNode,
2514                          NewNode->StartingVpn,
2515                          NewNode->EndingVpn);

2516

2517          if (PrevNode != NULL) {
2518              PRINT ("\tPrevNode %p %x %x\n",
2519                          PrevNode,
2520                          PrevNode->StartingVpn,
2521                          PrevNode->EndingVpn);
2522          }

2523

2524          if (NextNode != NULL) {
2525              PRINT ("\tNextNode %p %x %x\n",
2526                          NextNode,
2527                          NextNode->StartingVpn,
2528                          NextNode->EndingVpn);
2529          }

2530

2531      } while (TRUE);

2532
```

```
2533        PRINT ("NumberGenericTableElements␣=␣0x%x,␣Depth␣=␣0x%x\n",
2534            Table->NumberGenericTableElements,
2535            Table->DepthOfTree);
2536
2537        return;
2538 }
2539 #endif
2540
2541 #if defined (_USERMODE)
2542
2543 MMADDRESS_NODE MiBalancedLinks;
2544
2545 MM_AVL_TABLE MiAvlTable;
2546 MM_AVL_TABLE MmSectionBasedRoot;
2547
2548 ULONG DeleteRandom = 1;
2549
2550 #if RANDOM
2551 #define NUMBER_OF_VADS 32
2552 #else
2553 #define NUMBER_OF_VADS 4
2554 #endif
2555
2556 int __cdecl
2557 main(
2558 int argc,
2559 PCHAR    argv[]
2560 )
2561 {
2562      ULONG i;
2563      PVOID StartingAddress;
2564      PVOID EndingAddress;
2565      NTSTATUS Status;
2566      PMMADDRESS_NODE NewNode;
2567 #if RANDOM
2568      PMMADDRESS_NODE PrevNode;
2569      ULONG RandomNumber = 0x99887766;
2570      ULONG_PTR DeleteVpn = 0;
2571 #endif
```

```
2572      PMM_AVL_TABLE Table;
2573      SIZE_T CapturedRegionSize;
2574
2575      UNREFERENCED_PARAMETER (argc);
2576      UNREFERENCED_PARAMETER (argv);
2577
2578  #if RANDOM
2579      Table = &MiAvlTable;
2580  #else
2581      Table = &MmSectionBasedRoot;
2582  #endif
2583
2584      MiInitializeVadTableAvl (Table);
2585
2586      for (i = 0; i < NUMBER_OF_VADS; i += 1) {
2587          NewNode = malloc (sizeof (MMADDRESS_NODE));
2588          ASSERT (((ULONG_PTR)NewNode & 0x3) == 0);
2589
2590          if (NewNode == NULL) {
2591              PRINT ("Malloc failed\n");
2592              exit (1);
2593          }
2594
2595          NewNode->u1.Parent = NULL;
2596          NewNode->LeftChild = NULL;
2597          NewNode->RightChild = NULL;
2598          NewNode->u1.Balance = 0;
2599
2600  #if RANDOM
2601          RandomNumber = RtlRandom (&RandomNumber);
2602
2603          CapturedRegionSize = (SIZE_T) (RandomNumber & 0x1FFFFF);
2604
2605          Status = MiFindEmptyAddressRangeInTree (CapturedRegionSize,
2606                                                  64 * 1024,        //
                                                        align
2607                                                  Table,
2608                                                  &PrevNode,
2609                                                  &StartingAddress);
```

```
2610
2611 #else
2612         CapturedRegionSize = 0x800000;
2613
2614         Status = MiFindEmptyAddressRangeDownBasedTree (
                 CapturedRegionSize ,
2615                                               (PVOID) 0x7f7effff ,
                                                          // highest
                                                      addr
2616                                               64 * 1024 ,         //
                                                      align
2617                                               Table ,
2618                                               &StartingAddress );
2619 #endif
2620
2621         if (!NT_SUCCESS (Status)) {
2622             PRINT ("Could not find empty addr range in tree for size
                     %p\n", CapturedRegionSize );
2623             free (NewNode );
2624             continue ;
2625         }
2626
2627 #if RANDOM
2628         EndingAddress = (PVOID)(((ULONG_PTR)StartingAddress +
2629                             CapturedRegionSize − 1L) | (PAGE_SIZE −
                                 1L));
2630 #else
2631         EndingAddress = (PVOID)(((ULONG_PTR)StartingAddress +
2632                             CapturedRegionSize − 1L));
2633 #endif
2634
2635         printf ("Inserting addr range in tree @ %p %p\n",
                 StartingAddress , EndingAddress );
2636
2637 #if RANDOM
2638         NewNode−>StartingVpn = MI_VA_TO_VPN (StartingAddress );
2639         NewNode−>EndingVpn = MI_VA_TO_VPN (EndingAddress );
2640 #else
2641         NewNode−>StartingVpn = (ULONG_PTR) StartingAddress ;
```

```
2642            NewNode->EndingVpn = (ULONG_PTR) EndingAddress;
2643  #endif
2644
2645            MiInsertNode (NewNode, Table);
2646
2647  #if RANDOM
2648            RandomNumber = RtlRandom (&RandomNumber);
2649
2650            if (RandomNumber & 0x3) {
2651                DeleteVpn = NewNode->StartingVpn;
2652            }
2653
2654            if (DeleteRandom && ((i & 0x3) == 0)) {
2655                NewNode = MiLocateAddressInTree (DeleteVpn, Table);
2656                printf ("Located␣node␣for␣random␣deletion␣-␣vpn␣%p␣@␣%p\n
                        ", DeleteVpn, NewNode);
2657
2658                if (NewNode != NULL) {
2659                    MiRemoveNode (NewNode, Table);
2660                    printf ("Removed␣random␣node␣for␣vpn␣%p␣@␣%p␣%p␣%p\n"
                                ,
2661                            DeleteVpn, NewNode, NewNode->StartingVpn, NewNode
                                ->EndingVpn);
2662                }
2663            }
2664  #endif
2665            printf ("\n");
2666        }
2667
2668        MiNodeTreeWalk (Table);
2669
2670        NewNode = MiLocateAddressInTree (5, Table);
2671        printf ("Located␣node␣for␣vpn␣5␣@␣%p\n", NewNode);
2672
2673        if (NewNode != NULL) {
2674            MiRemoveNode (NewNode, Table);
2675            printf ("Removed␣node␣for␣vpn␣5␣@␣%p\n", NewNode);
2676        }
2677
```

```
2678        NewNode = MiLocateAddressInTree (5, Table);
2679        printf("Located␣node␣for␣vpn␣5␣@␣%p\n", NewNode);
2680
2681        printf("all␣done,␣balmin=%x,␣balmax=%x\n", BalMin, BalMax);
2682
2683        return  0;
2684    }
2685
2686    #endif
```