

File: ./init/main.c

```
1 #include <physical.h>
2 #include <virtual.h>
3 #include <heap.h>
4 #include <cpu.h>
5 #include <log.h>
6 #include <debug.h>
7 #include <assert.h>
8 #include <timer.h>
9 #include <irq.h>
10 #include <thread.h>
11 #include <panic.h>
12 #include <stdlib.h>
13 #include <process.h>
14 #include <progload.h>
15 #include <dev.h>
16 #include <vfs.h>
17 #include <diskcache.h>
18 #include <transfer.h>
19 #include <fcntl.h>
20 #include <console.h>
21 #include <swapfile.h>
22 #include <diskutil.h>
23 #include <string.h>
24 #include <filesystem.h>
25 #include <driver.h>
26
27 /*
28  * Next steps:
29  * - program loader / dynamic linker
30  * - system call interface (KRNLIB.LIB)
31  * - C standard library
32  * - complete-enough CLI OS
33  *   terminal that supports pipes, redirection and background processes
34  *   cd, ls/dir, type, mkdir, rm, more, rename, copy, tree, mkfifo, pause, rmtree, rmdir, cls, copytree, link,
35  *   ...ttyname, sleep, exit
36  *   port zlib, nasm
37  * - floppy driver
38  * - FAT32 driver
39  * - floating point support (init and task switching)
40  * - disk caching
41  * - shutdown needs to close the entire VFS tree (e.g. so buffers can be flushed, etc).
42  * - recycling vnodes if opening same file more than once
43  * - initrd and boot system
44  * - more syscalls
45  * - document exactly what conditions need to be checked in the vnode_ops layer, and which ones are taken care of by
46  *   the VFS layer, so we don't get people checking the same thing twice
47  * - check all E... return codes...
48  * - VnodeOpWait, select/poll syscalls
49  * - everyone create vnodes and open files willy-nilly - check the reference counting, especially on closing is
50  *   all correct (especially around the virtual memory manager...). does CloseFile do what you expect??
51  *   I THINK OPEN FILES SHOULD HOLD INCREMENT THE VNODE REFERENCE ON CREATION, AND DECREMENT WHEN THE OPEN FILE
52  *   GOES TO ZERO.
53  *
54  *   that way it will go a bit like this:
55  *
56  *       after call to:      vnode refs      openfile refs
57  *       Create vnode        1                  0
58  *       Create openfile     2                  1
59  *       CloseFile()         1                  0 -> gets destroyed
60  *       ->                  0
61  *
62  *       CloseFile() closes both the openfile and the vnode.
63  *
64  *   And then if the VMM gets involved...
65  *
66  *       after call to:      vnode refs      openfile refs
67  *       Create vnode        1                  0
68  *       Create openfile     2                  1
69  *       MapVirt(10 pgs)     2                  11
70  *       CloseFile()         1                  10
71  *       UnmapVirt(10 pgs)   1                  0 -> gets destroyed, and in doing so,
72  *       ...                 -> 0 gets destroyed
73  *
74  *   OK - that's been implemented now... now to see if it works...
75  *
76  * - MAP_FIXED, MAP_SHARED
77  *
78  */
79
80 void DummyAppThread(void*) {
81     PutsConsole("drv0:> ");
82
83     struct open_file* con;
84     OpenFile("con:", O_RDONLY, 0, &con);
85
86     while (true) {
87         char bf[302];
88         inline_memset(bf, 0, 302);
89         struct transfer tr = CreateKernelTransfer(bf, 301, 0, TRANSFER_READ);
90         ReadFile(con, &tr);
91         PutsConsole("Command not found: ");
92         PutsConsole(bf);
93         PutsConsole("\n");
94
95         if (bf[0] == 'u' || bf[0] == 'U') {
96             CreateUsermodeProcess(NULL, "sys:/init.exe");
97
98         } else if (bf[0] == 'p' || bf[0] == 'P') {
99             Panic(PANIC_MANUALLY_INITIATED);
100
101         } else if (bf[0] == 'e' || bf[0] == 'E') {
102             MapVirt(0, 0, 8 * 4096, VM_LOCK | VM_READ, NULL, 0);
103         }
104
105         PutsConsole("drv0:> ");
106     }
107 }
108
109 void InitUserspace(void) {
110     size_t free = GetFreePhysKilobytes();
111     size_t total = GetTotalPhysKilobytes();
112     DbgScreenPrintf("NOS Kernel\nCopyright Alex Boxall 2022-2023\n\n%d / %d KB used (%d%% free)\n\n", total - free, total, 100 * (free) / total);
113     CreateThread(DummyAppThread, NULL, GetVas(), "dummy app");
114 }
115
116 void InitThread(void*) {
117     ReinitHeap();
118     InitRandomDevice();
119     InitNullDevice();
120     InitConsole();
121     InitProcess();
122     InitDiskCaches();
123
124     InitFilesystemTable();
125     ArchInitDev(false);
126
127     struct open_file* sys_folder;
128     int res = OpenFile("drv0:/System", O_RDONLY, 0, &sys_folder);
```

```

129     if (res != 0) {
130         PanicEx (PANIC_NO_FILESYSTEM, "sys A");
131     }
132     res = AddVfsMount (sys_folder->node, "sys");
133     if (res != 0) {
134         PanicEx (PANIC_NO_FILESYSTEM, "sys B");
135     }
136
137     struct open_file* swapfile;
138     res = OpenFile ("raw-hd0:/part1", O_RDWR, 0, &swapfile);
139     if (res != 0) {
140         PanicEx (PANIC_NO_FILESYSTEM, "swapfile A");
141     }
142     res = AddVfsMount (swapfile->node, "swap");
143     if (res != 0) {
144         PanicEx (PANIC_NO_FILESYSTEM, "swapfile B");
145     }
146
147     InitSwapfile();
148     InitSymbolTable();
149     ArchInitDev (true);
150     InitProgramLoader();
151     InitUserspace();
152
153     MarkTfwStartPoint (TFW_SP_ALL_CLEAR);
154
155     while (true) {
156         /*
157          * We crash in strange and rare conditions if this thread's stack gets removed, so we will
158          * ensure we don't terminate it.
159          */
160         SleepMilli (100000);
161     }
162 }
163
164 #include <machine/portio.h>
165 static void InitSerialDebugging(void) {
166     const int PORT = 0x3F8;
167     outb (PORT + 1, 0x00); // Disable all interrupts
168     outb (PORT + 3, 0x80); // Enable DLAB (set baud rate divisor)
169     outb (PORT + 0, 0x03); // Set divisor to 3 (lo byte) 38400 baud
170     outb (PORT + 1, 0x00); // (hi byte)
171     outb (PORT + 3, 0x03); // 8 bits, no parity, one stop bit
172     outb (PORT + 2, 0xC7); // Enable FIFO, clear them, with 14-byte threshold
173     outb (PORT + 4, 0x0B); // IRQs enabled, RTS/DSR set
174     outb (PORT + 4, 0x1E); // Set in loopback mode, test the serial chip
175     outb (PORT + 0, 0xAE); // Test serial chip (send byte 0xAE and check if serial returns same byte)
176
177     // Check if serial is faulty (i.e: not same byte as sent)
178     if (inb (PORT + 0) != 0xAE) {
179         return;
180     }
181
182     // If serial is not faulty set it in normal operation mode
183     // (not-loopback with IRQs enabled and OUT#1 and OUT#2 bits enabled)
184     outb (PORT + 4, 0x0F);
185 }
186
187 void KernelMain(void) {
188     InitSerialDebugging();
189
190     LogWriteSerial ("KernelMain: kernel is initialising...\n");
191
192     /*
193      * Allows us to call GetCpu(), which allows IRQL code to work. Anything which uses
194      * IRQL (i.e. the whole system) relies on this, so this must be done first.
195      */
196     InitCpuTable();
197     assert (GetIrql() == IRQL_STANDARD);
198
199     /*
200      * Initialise the testing framework if we're in debug mode.
201      */
202     InitTfw();
203     MarkTfwStartPoint (TFW_SP_INITIAL);
204
205     InitPhys();
206     MarkTfwStartPoint (TFW_SP_AFTER_PHYS);
207
208     /*
209      * Allows deferments of functions to actually happen. IRQL is still usable beforehand though.
210      */
211     InitIrql();
212     InitVfs();
213     InitTimer();
214     InitScheduler();
215     InitDiskUtil();
216
217     InitHeap();
218     MarkTfwStartPoint (TFW_SP_AFTER_HEAP);
219
220     InitBootstrapCpu();
221     MarkTfwStartPoint (TFW_SP_AFTER_BOOTSTRAP_CPU);
222
223     InitVirt();
224     MarkTfwStartPoint (TFW_SP_AFTER_VIRT);
225
226     ReinitPhys();
227     MarkTfwStartPoint (TFW_SP_AFTER_PHYS_REINIT);
228
229     InitOtherCpu();
230     MarkTfwStartPoint (TFW_SP_AFTER_ALL_CPU);
231
232     CreateThreadEx (InitThread, NULL, GetVas(), "init", NULL, SCHEDULE_POLICY_FIXED, FIXED_PRIORITY_KERNEL_NORMAL, 0);
233     StartMultitasking();
234 }

```

File: `Jad/tthreadlist.c`

```

1  #include <common.h>
2  #include <threadlist.h>
3  #include <heap.h>
4  #include <assert.h>
5  #include <thread.h>
6  #include <panic.h>
7  #include <log.h>
8  #include <string.h>
9
10 void ThreadListInit(struct thread_list* list, int index) {
11     inline memset(list, 0, sizeof(struct thread_list));
12     list->index = index;
13 }
14
15 void ThreadListInsert(struct thread_list* list, struct thread* thread) {
16     #ifndef NDEBUG
17         if (ThreadListContains(list, thread)) {
18             assert(!ThreadListContains(list, thread));
19         }
20     #endif
21
22     if (list->tail == NULL) {
23         assert(list->head == NULL);
24         list->head = thread;
25     } else {
26         list->tail->next[list->index] = thread;
27     }
28
29     list->tail = thread;
30     thread->next[list->index] = NULL;
31 }
32
33 static int ThreadListGetIndex(struct thread_list* list, struct thread* thread) {
34     struct thread* iter = list->head;
35     int i = 0;
36     while (iter != NULL) {
37         if (iter == thread) {
38             return i;
39         }
40         ++i;
41         iter = iter->next[list->index];
42     }
43     return -1;
44 }
45
46 bool ThreadListContains(struct thread_list* list, struct thread* thread) {
47     return ThreadListGetIndex(list, thread) != -1;
48 }
49
50 static void ProperDelete(struct thread_list* list, struct thread* iter, struct thread* prev) {
51     if (iter == list->head) {
52         list->head = list->head->next[list->index];
53     } else {
54         prev->next[list->index] = iter->next[list->index];
55     }
56
57     if (iter == list->tail) {
58         list->tail = prev;
59     }
60
61     if (list->head == NULL) {
62         assert(list->tail == NULL);
63     }
64 }
65
66 static void ThreadListDeleteIndex(struct thread_list* list, int index) {
67     struct thread* iter = list->head;
68     struct thread* prev = NULL;
69
70     assert(index >= 0);
71
72     int i = 0;
73     while (iter != NULL) {
74         if (i == index) {
75             ProperDelete(list, iter, prev);
76             return;
77         }
78         ++i;
79         prev = iter;
80         iter = iter->next[list->index];
81     }
82
83     Panic(PANIC_THREAD_LIST);
84 }
85
86 struct thread* ThreadListDeleteTop(struct thread_list* list) {
87     struct thread* top = list->head;
88     if (list->head == list->tail) {
89         list->tail = NULL;
90     }
91     list->head = list->head->next[list->index];
92
93     if (list->head == NULL) {
94         assert(list->tail == NULL);
95     }
96
97     return top;
98 }
99
100 void ThreadListDelete(struct thread_list* list, struct thread* thread) {
101     ThreadListDeleteIndex(list, ThreadListGetIndex(list, thread));
102 }
103

```

File: /adt/avl.c

```

1  #include <common.h>
2  #include <assert.h>
3  #include <heap.h>
4  #include <avl.h>
5  #include <log.h>
6
7  struct avl_node {
8     struct avl_node* left;
9     struct avl_node* right;
10     void* data;
11 };
12
13 struct avl_tree {
14     int size;
15     struct avl_node* root;
16     avl_deletion_handler deletion_handler;
17     avl_comparator equality_handler;
18 };
19
20 static struct avl_node* AvlCreateNode(void* data, struct avl_node* restrict left, struct avl_node* restrict right) {
21     struct avl_node* tree = AllocHeap(sizeof(struct avl_node));

```

```

22     tree->left = left;
23     tree->right = right;
24     tree->data = data;
25     return tree;
26 }
27
28 static int AvlGetHeight(struct avl_node* tree) {
29     if (tree == NULL) {
30         return 0;
31     }
32
33     return 1 + MAX(AvlGetHeight(tree->left), AvlGetHeight(tree->right));
34 }
35
36 static int AvlGetBalance(struct avl_node* tree) {
37     if (tree == NULL) {
38         return 0;
39     }
40
41     return AvlGetHeight(tree->left) - AvlGetHeight(tree->right);
42 }
43
44 static struct avl_node* AvlRotateLeft(struct avl_node* tree) {
45     struct avl_node* new_root = tree->right;
46     struct avl_node* new_right = new_root->left;
47     new_root->left = tree;
48     tree->right = new_right;
49     return new_root;
50 }
51
52 static struct avl_node* AvlRotateRight(struct avl_node* tree) {
53     struct avl_node* new_root = tree->left;
54     struct avl_node* new_left = new_root->right;
55     new_root->right = tree;
56     tree->left = new_left;
57     return new_root;
58 }
59
60 static struct avl_node* AvlBalance(struct avl_node* tree) {
61     if (tree == NULL) {
62         return NULL;
63     }
64
65     int bf = AvlGetBalance(tree);
66     assert(bf >= -2 && bf <= 2);
67
68     if (bf == -2) {
69         if (AvlGetBalance(tree->right) == 1) {
70             tree->right = AvlRotateRight(tree->right);
71         }
72         return AvlRotateLeft(tree);
73     }
74
75     if (bf == 2) {
76         if (AvlGetBalance(tree->left) == -1) {
77             tree->left = AvlRotateLeft(tree->left);
78         }
79         return AvlRotateRight(tree);
80     }
81     return tree;
82 }
83
84 static struct avl_node* AvlInsert(struct avl_node* tree, void* data, avl_comparator comparator) {
85     struct avl_node* new_tree;
86
87     // TODO: surely there's a better way that involves less node creation and
88     // deletion...
89
90     assert(comparator != NULL);
91     assert(tree != NULL);
92
93     if (comparator(data, tree->data) < 0) {
94         struct avl_node* left_tree;
95         if (tree->left == NULL) {
96             left_tree = AvlCreateNode(data, NULL, NULL);
97         } else {
98             left_tree = AvlInsert(tree->left, data, comparator);
99         }
100
101         new_tree = AvlCreateNode(tree->data, left_tree, tree->right);
102     } else {
103         struct avl_node* right_tree;
104         if (tree->right == NULL) {
105             right_tree = AvlCreateNode(data, NULL, NULL);
106         } else {
107             right_tree = AvlInsert(tree->right, data, comparator);
108         }
109
110         new_tree = AvlCreateNode(tree->data, tree->left, right_tree);
111     }
112
113     FreeHeap(tree);
114
115     return AvlBalance(new_tree);
116 }
117
118 static struct avl_node* AvlDelete(struct avl_node* tree, void* data, avl_comparator comparator) {
119     if (tree == NULL) {
120         return NULL;
121     }
122
123     struct avl_node* to_free = NULL;
124
125     if (comparator(data, tree->data) < 0) {
126         tree->left = AvlDelete(tree->left, data, comparator);
127     } else if (comparator(data, tree->data) > 0) {
128         tree->right = AvlDelete(tree->right, data, comparator);
129     } else if (tree->left == NULL) {
130         to_free = tree;
131         tree = tree->right;
132     } else if (tree->right == NULL) {
133         to_free = tree;
134         tree = tree->left;
135     } else {
136         struct avl_node* node = tree->right;
137         while (node->left != NULL) {
138             node = node->left;
139         }
140
141         tree->data = node->data;
142         tree->right = AvlDelete(tree->right, node->data, comparator);
143     }
144
145     /*
146
147
148
149
150
151

```

```

152     * If NULL is passed in to FreeHeap, nothing happens (which is what we want).
153     */
154     FreeHeap(to_free);
155
156     return AvlBalance(tree);
157 }
158
159 /**
160  * Given an object, find it in the AVL tree and return it. This is useful if the comparator only compares
161  * part of the object, and so the entire object can be retrieved by searching for only part of it.
162  */
163 static void* AvlGet(struct avl_node* tree, void* data, avl_comparator comparator) {
164     if (tree == NULL) {
165         return NULL;
166     }
167
168     if (comparator(tree->data, data) == 0) {
169         /*
170          * Must return 'tree->data', (and not 'data'), as tree->data != data if there is a custom comparator.
171          */
172         return tree->data;
173     }
174
175     void* left = AvlGet(tree->left, data, comparator);
176     if (left != NULL) {
177         return left;
178     }
179     return AvlGet(tree->right, data, comparator);
180 }
181
182 static void AvlPrint(struct avl_node* tree, void(*printer)(void*)) {
183     if (tree == NULL) {
184         return;
185     }
186     AvlPrint(tree->left, printer);
187     if (printer == NULL) {
188         LogWriteSerial("[0x%X]", "\n", tree->data);
189     } else {
190         printer(tree->data);
191     }
192     AvlPrint(tree->right, printer);
193 }
194
195 static bool AvlContains(struct avl_node* tree, void* data, avl_comparator comparator) {
196     if (tree == NULL) {
197         return false;
198     }
199
200     if (comparator(tree->data, data) == 0) {
201         return true;
202     }
203
204     return AvlContains(tree->left, data, comparator) || AvlContains(tree->right, data, comparator);
205 }
206
207 static void AvlDestroy(struct avl_node* tree, avl_deletion_handler handler) {
208     if (tree == NULL) {
209         return;
210     }
211
212     AvlDestroy(tree->left, handler);
213     AvlDestroy(tree->right, handler);
214     if (handler != NULL) {
215         handler(tree->data);
216     }
217     FreeHeap(tree);
218 }
219
220 static int AvlDefaultComparator(void* a, void* b) {
221     if (a == b) return 0;
222     return (a < b) ? -1 : 1;
223 }
224
225 struct avl_tree* AvlTreeCreate(void) {
226     struct avl_tree* tree = AllocHeap(sizeof(struct avl_tree));
227     tree->size = 0;
228     tree->root = NULL;
229     tree->deletion_handler = NULL;
230     tree->equality_handler = AvlDefaultComparator;
231     return tree;
232 }
233
234 avl_deletion_handler AvlTreeSetDeletionHandler(struct avl_tree* tree, avl_deletion_handler handler) {
235     avl_deletion_handler ret = tree->deletion_handler;
236     tree->deletion_handler = handler;
237     return ret;
238 }
239
240 avl_comparator AvlTreeSetComparator(struct avl_tree* tree, avl_comparator handler) {
241     avl_comparator ret = tree->equality_handler;
242     tree->equality_handler = handler;
243     return ret;
244 }
245
246 void AvlTreeInsert(struct avl_tree* tree, void* data) {
247     if (tree->root == NULL) {
248         tree->root = AvlCreateNode(data, NULL, NULL);
249     } else {
250         tree->root = AvlInsert(tree->root, data, tree->equality_handler);
251     }
252     tree->size++;
253 }
254
255 void AvlTreeDelete(struct avl_tree* tree, void* data) {
256     tree->root = AvlDelete(tree->root, data, tree->equality_handler);
257     tree->size--;
258 }
259
260 bool AvlTreeContains(struct avl_tree* tree, void* data) {
261     return AvlContains(tree->root, data, tree->equality_handler);
262 }
263
264 void* AvlTreeGet(struct avl_tree* tree, void* data) {
265     return AvlGet(tree->root, data, tree->equality_handler);
266 }
267
268 int AvlTreeSize(struct avl_tree* tree) {
269     return tree->size;
270 }
271
272 void AvlTreeDestroy(struct avl_tree* tree) {
273     AvlDestroy(tree->root, tree->deletion_handler);
274     FreeHeap(tree);
275 }
276
277 struct avl_node* AvlTreeGetRootNode(struct avl_tree* tree) {
278     return tree->root;
279 }
280
281 struct avl_node* AvlTreeGetLeft(struct avl_node* node) {

```

```

282     if (node == NULL) {
283         return NULL;
284     }
285     return node->left;
286 }
287
288
289 struct avl_node* AvlTreeGetRight(struct avl_node* node) {
290     if (node == NULL) {
291         return NULL;
292     }
293     return node->right;
294 }
295
296
297 void* AvlTreeGetData(struct avl_node* node) {
298     if (node == NULL) {
299         return NULL;
300     }
301     return node->data;
302 }
303
304
305 void AvlTreePrint(struct avl_tree* tree, void(*printer)(void*)) {
306     AvlPrint(tree->root, printer);
307 }

```

File: /adt/priorityqueue.c

```

1
2
3 #include <heap.h>
4 #include <string.h>
5 #include <log.h>
6 #include <panic.h>
7 #include <assert.h>
8 #include <priorityqueue.h>
9
10 /*
11  * Implements the max-heap and min-heap data structures. To avoid confusion with the
12  * heap memory manager, it is referred to as a priority queue.
13  */
14
15 struct priority_queue {
16     int capacity;
17     int size;
18     int element_width;
19     int qwords_per_element; // includes the + 1 for the priority
20     bool max;
21     uint64_t* array; // length is: capacity * qwords_per_element
22 };
23
24 struct priority_queue* PriorityQueueCreate(int capacity, bool max, int element_width) {
25     assert(capacity > 0);
26     assert(element_width > 0);
27
28     struct priority_queue* queue = AllocHeap(sizeof(struct priority_queue));
29     queue->capacity = capacity;
30     queue->size = 0;
31     queue->element_width = element_width;
32     queue->qwords_per_element = 1 + (element_width + sizeof(uint64_t) - 1) / sizeof(uint64_t);
33     queue->max = max;
34     queue->array = AllocHeap(sizeof(uint64_t) * queue->qwords_per_element * capacity);
35     return queue;
36 }
37
38 void PriorityQueueDestroy(struct priority_queue* queue) {
39     FreeHeap(queue->array);
40     FreeHeap(queue);
41 }
42
43 static void SwapElements(struct priority_queue* queue, int a, int b) {
44     a *= queue->qwords_per_element;
45     b *= queue->qwords_per_element;
46
47     for (int i = 0; i < queue->qwords_per_element; ++i) {
48         uint64_t tmp = queue->array[a];
49         queue->array[a] = queue->array[b];
50         queue->array[b] = tmp;
51         ++a; ++b;
52     }
53 }
54
55 static void Heapify(struct priority_queue* queue, int i) {
56     int extreme = i;
57     int left = i * 2 + 1;
58     int right = left + 1;
59
60     if (left < queue->size) {
61         if ((queue->max && queue->array[left * queue->qwords_per_element] > queue->array[extreme * queue->qwords_per_element]) || (!queue->max && queue->array[left * queue->qwords_per_element] < queue->array[extreme * queue->qwords_per_element])) {
62             extreme = left;
63         }
64     }
65
66     if (right < queue->size) {
67         if ((queue->max && queue->array[right * queue->qwords_per_element] > queue->array[extreme * queue->qwords_per_element]) || (!queue->max && queue->array[right * queue->qwords_per_element] < queue->array[extreme * queue->qwords_per_element])) {
68             extreme = right;
69         }
70     }
71
72     if (i != extreme) {
73         SwapElements(queue, i, extreme);
74         Heapify(queue, extreme);
75     }
76 }
77
78 void PriorityQueueInsert(struct priority_queue* queue, void* elem, uint64_t priority) {
79     if (queue->size == queue->capacity) {
80         // I think this can happen when the OS is running too slowly! (the deferred function buffer actually fills up
81         // and overflows). Testing on real H/W from 1996, debug serial port accesses took about a second per character!!
82         PanicEx(PANIC_PRIORITY_QUEUE, "insert called when full");
83     }
84
85     int i = queue->size++;
86     queue->array[i * queue->qwords_per_element] = priority;
87     inline_memcpy(queue->array + i * queue->qwords_per_element + 1, elem, queue->element_width);
88
89     if (queue->max) {
90         while (i != 0 && queue->array[(i - 1) / 2] * queue->qwords_per_element < queue->array[i * queue->qwords_per_element]) {
91             SwapElements(queue, (i - 1) / 2, i);
92             i = (i - 1) / 2;
93         }
94     } else {
95         while (i != 0 && queue->array[(i - 1) / 2] * queue->qwords_per_element > queue->array[i * queue->qwords_per_element]) {
96             SwapElements(queue, (i - 1) / 2, i);
97             i = (i - 1) / 2;
98         }
99     }
100 }

```

```

98 }
99
100 /*
101  * returns the priority, and a REFERENCE to the data IN THE PRIORITY QUEUE. It is NOT a copy!
102  * This is why Pop() can't return it, because popping it overwrites the data!
103  */
104 struct priority_queue_result PriorityQueuePeek(struct priority_queue* queue) {
105     if (queue->size == 0)
106         PanicEx(PANIC_PRIORITY_QUEUE, "peek called on empty");
107 }
108
109 struct priority_queue_result retv;
110 retv.priority = queue->array[0];
111 retv.data = (void*) (queue->array + 1);
112 return retv;
113 }
114
115 /*
116  * doesn't return the value, as the value would have been erased in the pop (PriorityQueue doesn't
117  * allocate memory, as it needs to be used in high IRQ/L situations).
118  */
119 void PriorityQueuePop(struct priority_queue* queue) {
120     if (queue->size == 0)
121         PanicEx(PANIC_PRIORITY_QUEUE, "pop called on empty");
122 }
123
124 --queue->size;
125 for (int i = 0; i < queue->qwords_per_element; ++i) {
126     queue->array[i] = queue->array[queue->size * queue->qwords_per_element + i];
127 }
128
129 Heapify(queue, 0);
130 }
131
132 int PriorityQueueGetCapacity(struct priority_queue* queue) {
133     return queue->capacity;
134 }
135
136 int PriorityQueueGetUsedSize(struct priority_queue* queue) {
137     return queue->size;
138 }

```

File: /adt/stackadt.c

```

1 #include <common.h>
2 #include <linkedlist.h>
3 #include <stackadt.h>
4 #include <heap.h>
5 #include <assert.h>
6 #include <panic.h>
7
8 struct stack_adt {
9     struct linked_list* list;
10 };
11
12 struct stack_adt* StackAdtCreate(void) {
13     struct stack_adt* stack = AllocHeap(sizeof(struct stack_adt));
14     stack->list = LinkedListCreate();
15     return stack;
16 }
17
18 void StackAdtDestroy(struct stack_adt* stack) {
19     LinkedListDestroy(stack->list);
20     FreeHeap(stack);
21 }
22
23 void StackAdtPush(struct stack_adt* stack, void* data) {
24     LinkedListInsertStart(stack->list, data);
25 }
26
27 void* StackAdtPeek(struct stack_adt* stack) {
28     return LinkedListGetDataFromNode(LinkedListGetFirstNode(stack->list));
29 }
30
31 void* StackAdtPop(struct stack_adt* stack) {
32     void* data = StackAdtPeek(stack);
33     LinkedListDeleteIndex(stack->list, 0);
34     return data;
35 }
36
37 int StackAdtSize(struct stack_adt* stack) {
38     return LinkedListSize(stack->list);
39 }

```

File: /adt/blockingbuffer.c

```

1
2 #include <heap.h>
3 #include <assert.h>
4 #include <common.h>
5 #include <spinlock.h>
6 #include <semaphore.h>
7 #include <errno.h>
8 #include <thread.h>
9 #include <panic.h>
10 #include <log.h>
11 #include <irq.h>
12
13 struct blocking_buffer {
14     uint8_t* buffer;
15     int total_size;
16     int used_size;
17     int start_pos;
18     int end_pos;
19     struct semaphore* sem;
20     struct semaphore* reverse_sem;
21     struct spinlock lock;
22 };
23
24 struct blocking_buffer* BlockingBufferCreate(int size) {
25     assert(size > 0);
26
27     struct blocking_buffer* buffer = AllocHeap(sizeof(struct blocking_buffer));
28     buffer->buffer = AllocHeap(size);
29     buffer->total_size = size;
30     buffer->used_size = 0;
31     buffer->start_pos = 0;
32     buffer->end_pos = 0;
33     buffer->sem = CreateSemaphore("bb get", size, size);
34     buffer->reverse_sem = CreateSemaphore("bb add", size, 0);
35     InitSpinlock(&buffer->lock, "blocking buffer", IRQL_SCHEDULER);
36     return buffer;
37 }
38
39 void BlockingBufferDestroy(struct blocking_buffer* buffer) {
40     FreeHeap(buffer->sem);
41     FreeHeap(buffer->buffer);
42     FreeHeap(buffer);
43 }
44
45 int BlockingBufferAdd(struct blocking_buffer* buffer, uint8_t c, bool block) {
46     int res = AcquireSemaphore(buffer->reverse_sem, block ? -1 : 0);
47
48     if (!block && res != 0) {
49         return ENOBUFS;
50     }
51
52     AcquireSpinlockIrql(&buffer->lock);
53
54     assert(buffer->used_size != buffer->total_size);
55
56     buffer->buffer[buffer->end_pos] = c;
57     buffer->end_pos = (buffer->end_pos + 1) % buffer->total_size;
58     buffer->used_size++;
59
60     /*
61      * Wake up someone waiting for a character to enter the buffer - or make it so
62      * next time someone wants a character they can grab it straight away.
63      */
64     ReleaseSpinlockIrql(&buffer->lock);
65     ReleaseSemaphore(buffer->sem);
66     return 0;
67 }
68
69 static uint8_t BlockingBufferGetAfterAcquisition(struct blocking_buffer* buffer) {
70     AcquireSpinlockIrql(&buffer->lock);
71
72     uint8_t c = buffer->buffer[buffer->start_pos];
73     buffer->start_pos = (buffer->start_pos + 1) % buffer->total_size;
74     buffer->used_size--;
75
76     ReleaseSpinlockIrql(&buffer->lock);
77     ReleaseSemaphore(buffer->reverse_sem);
78
79     return c;
80 }
81
82 uint8_t BlockingBufferGet(struct blocking_buffer* buffer) {
83     /*
84      * Wait for there to be something to actually read.
85      */
86     AcquireSemaphore(buffer->sem, -1);
87     return BlockingBufferGetAfterAcquisition(buffer);
88 }
89
90 int BlockingBufferTryGet(struct blocking_buffer* buffer, uint8_t* c) {
91     assert(c != NULL);
92
93     int result = AcquireSemaphore(buffer->sem, 0);
94     if (result == 0) {
95         *c = BlockingBufferGetAfterAcquisition(buffer);
96         return 0;
97     } else {
98         return result;
99     }
100 }
101

```

File: /adt/linkedlist.c

```

1 #include <common.h>
2 #include <linkedlist.h>
3 #include <heap.h>
4 #include <assert.h>
5 #include <panic.h>
6
7 struct linked_list_node {
8     void* data;
9     struct linked_list_node* next;
10 };
11
12 struct linked_list {
13     int size;
14     struct linked_list_node* head;
15     struct linked_list_node* tail;
16 };
17
18 struct linked_list* LinkedListCreate(void) {
19     struct linked_list* list = AllocHeap(sizeof(struct linked_list));
20     list->size = 0;
21     list->head = NULL;
22     list->tail = NULL;

```



```

23     return list;
24 }
25
26 void LinkedListInsertStart(struct linked_list* list, void* data) {
27     struct linked_list_node* node = AllocHeap(sizeof(struct linked_list_node));
28     node->data = data;
29     node->next = list->tail;
30
31     if (list->head == NULL) {
32         assert(list->tail == NULL);
33         list->tail = node;
34     }
35
36     list->head = node;
37     list->size++;
38 }
39
40 void LinkedListInsertEnd(struct linked_list* list, void* data) {
41     if (list->tail == NULL) {
42         assert(list->head == NULL);
43         list->tail = AllocHeap(sizeof(struct linked_list_node));
44         list->head = list->tail;
45     }
46     else {
47         list->tail->next = AllocHeap(sizeof(struct linked_list_node));
48         list->tail = list->tail->next;
49     }
50
51     list->tail->data = data;
52     list->tail->next = NULL;
53     list->size++;
54 }
55
56 bool LinkedListContains(struct linked_list* list, void* data) {
57     return LinkedListGetIndex(list, data) != -1;
58 }
59
60 int LinkedListGetIndex(struct linked_list* list, void* data) {
61     struct linked_list_node* iter = list->head;
62     int i = 0;
63     while (iter != NULL) {
64         if (iter->data == data) {
65             return i;
66         }
67         ++i;
68         iter = iter->next;
69     }
70     return -1;
71 }
72
73 void* LinkedListGetData(struct linked_list* list, int index) {
74     struct linked_list_node* iter = list->head;
75     int i = 0;
76     while (iter != NULL) {
77         if (i == index) {
78             return iter->data;
79         }
80         ++i;
81         iter = iter->next;
82     }
83     Panic(PANIC_LINKED_LIST);
84 }
85
86 static void ProperDelete(struct linked_list* list, struct linked_list_node* iter, struct linked_list_node* prev) {
87     if (iter == list->head) {
88         list->head = list->head->next;
89     }
90     else {
91         prev->next = iter->next;
92     }
93
94     if (iter == list->tail) {
95         list->tail = prev;
96     }
97
98     FreeHeap(iter);
99     list->size--;
100 }
101
102 bool LinkedListDeleteIndex(struct linked_list* list, int index) {
103     if (index >= list->size || index < 0) {
104         return false;
105     }
106
107     struct linked_list_node* iter = list->head;
108     struct linked_list_node* prev = NULL;
109
110     int i = 0;
111     while (iter != NULL) {
112         if (i == index) {
113             ProperDelete(list, iter, prev);
114             return true;
115         }
116         ++i;
117         prev = iter;
118         iter = iter->next;
119     }
120
121     return false;
122 }
123
124 bool LinkedListDeleteData(struct linked_list* list, void* data) {
125     return LinkedListDeleteIndex(list, LinkedListGetIndex(list, data));
126 }
127
128 int LinkedListSize(struct linked_list* list) {
129     return list->size;
130 }
131
132 void LinkedListDestroy(struct linked_list* list) {
133     while (list->size > 0) {
134         LinkedListDeleteIndex(list, 0);
135     }
136     FreeHeap(list);
137 }
138
139 struct linked_list_node* LinkedListGetFirstNode(struct linked_list* list) {
140     if (list == NULL) {
141         Panic(PANIC_LINKED_LIST);
142     }
143     return list->head;
144 }
145
146 struct linked_list_node* LinkedListGetNextNode(struct linked_list_node* prev_node) {
147     if (prev_node == NULL) {
148         Panic(PANIC_LINKED_LIST);
149     }
150     return prev_node->next;
151 }
152
153 void* LinkedListGetDataFromNode(struct linked_list_node* node) {

```

```

153     if (node == NULL) {
154         Panic(PANIC_LINKED_LIST);
155     }
156     return node->data;
157 }

```

File: ./vfs/openfile.c

```

1  #include <openfile.h>
2  #include <spinlock.h>
3  #include <assert.h>
4  #include <heap.h>
5  #include <iqrl.h>
6  #include <vfs.h>
7
8  /**
9   * Creates a new open file from an opened vnode. An open file is used to link a vnode with corresponding data,
10   * about a particular instance of opening a file: such as a seek position, and ability to read or write; and is
11   * used to maintain file descriptor tables for the C userspace library.
12   *
13   * Open files maintain a reference counter that can be incremented and decremented with ReferenceOpenFile and
14   * DereferenceOpenFile. A newly created open file has a reference count of 1. When the count reaches 0, the memory
15   * is freed.
16   *
17   * @param node      The already-open vnode to wrap with additional data
18   * @param mode      The Unix-permissions passed to the OpenFile when the vnode was opened. Ignored by the kernel
19   *                  so far.
20   * @param flag      The flags that were passed to OpenFile when the vnode was opened. Should be a bitfield consisting of
21   *                  zero or more of: O_RDONLY, O_WRONLY, O_TRUNC, O_CREAT. See OpenFile for details. The storage of these
22   *                  flags in an open file should only be of interest to usermode programs - the flags here may be zero
23   *                  if the open file was created within the kernel, so the flags should be ignored within the kernel.
24   * @param can_read  Whether or not this open file is allowed to make read operations on the underlying vnode
25   * @param can_write Whether or not this open file is allowed to make write operations on the underlying vnode
26   *
27   * @return A pointer to the newly created open file.
28   */
29 struct open_file* CreateOpenFile(struct vnode* node, int mode, int flags, bool can_read, bool can_write) {
30     MAX_IQRL(IQRL_SCHEDULER);
31
32     struct open_file* file = AllocHeap(sizeof(struct open_file));
33     file->reference_count = 1;
34     file->node = node;
35     file->can_read = can_read;
36     file->can_write = can_write;
37     file->initial_mode = mode;
38     file->flags = flags;
39     file->seek_position = 0;
40     InitSpinlock(&file->reference_count_lock, "open file", IQRL_SCHEDULER);
41
42     ReferenceVnode(node);
43
44     return file;
45 }
46
47 /**
48   * Increments the reference counter for an opened file. This should be called everytime a reference to
49   * the open file is kept, so that its memory can be managed correctly.
50   *
51   * @param file The open file to reference
52   */
53 void ReferenceOpenFile(struct open_file* file) {
54     MAX_IQRL(IQRL_SCHEDULER);
55     assert(file != NULL);
56
57     AcquireSpinlockIrql(&file->reference_count_lock);
58     file->reference_count++;
59     ReleaseSpinlockIrql(&file->reference_count_lock);
60 }
61
62 /**
63   * Decrements the reference counter for an opened file. Should be called whenever a reference to the open
64   * file is removed. If the reference counter reaches zero, the memory behind the open file will be freed.
65   * The underlying vnode within the open file is not dereferenced - this should be done prior to calling this
66   * function.
67   *
68   * @param file The open file to dereference
69   */
70 void DereferenceOpenFile(struct open_file* file) {
71     MAX_IQRL(IQRL_SCHEDULER);
72     assert(file != NULL);
73
74     AcquireSpinlockIrql(&file->reference_count_lock);
75
76     assert(file->reference_count > 0);
77     file->reference_count--;
78
79     if (file->reference_count == 0) {
80         /*
81          * Must release the lock before we delete it so we can put interrupts back on
82          */
83         ReleaseSpinlockIrql(&file->reference_count_lock);
84         DereferenceVnode(file->node);
85         FreeHeap(file);
86         return;
87     }
88
89     ReleaseSpinlockIrql(&file->reference_count_lock);
90 }

```

File: ./vfs/diskutil.c

```

1  #include <diskutil.h>
2  #include <string.h>
3  #include <iqrl.h>
4  #include <assert.h>
5  #include <spinlock.h>
6  #include <vfs.h>
7  #include <xerrno.h>
8  #include <log.h>
9  #include <iqrl.h>
10 #include <partition.h>
11 #include <sys/stat.h>
12
13 /**
14   * Stores how many disks of each type have been allocated so far.
15   */
16 static int type_table[_DISKUTIL_NUM_TYPES];
17
18 /**
19   * Protects 'type_table' and 'next_mounted_disk_num'
20   */
21 static struct spinlock type_table_lock;
22

```

```

23 /*
24  * Filesystems get mounted to the VFS as drvX:, where X is an increasing number.
25  * This value stores the number the next disk gets.
26  */
27 static int next_mounted_disk_num = 0;
28
29 /*
30  * Maps a drive type to a string that will form part of the drive name.
31  */
32 static char* type_strings[_DISKUTIL_NUM_TYPES] = {
33     [DISKUTIL_TYPE_FIXED] = "hd",
34     [DISKUTIL_TYPE_FLOPPY] = "fd",
35     [DISKUTIL_TYPE_NETWORK] = "net",
36     [DISKUTIL_TYPE_OPTICAL] = "cd",
37     [DISKUTIL_TYPE_OTHER] = "other",
38     [DISKUTIL_TYPE_RAM] = "ram",
39     [DISKUTIL_TYPE_REMOVABLE] = "rm",
40     [DISKUTIL_TYPE_VIRTUAL] = "virt",
41 };
42
43 /**
44  * Initialises the disk utility functions. Must be called before any partitions
45  * are created or any drive names are generated.
46  */
47 void InitDiskUtil(void) {
48     EXACT_IRQL(IRQL_STANDARD);
49     InitSpinlock(&type_table_lock, "diskutil", IRQL_SCHEDULER);
50     memset(type_table, 0, sizeof(type_table));
51 }
52
53 /**
54  * Given a string and an integer less than 1000, it converts the integer to a
55  * string, and appends it to the end of the existing string, in place. The
56  * string should have enough buffer allocated to fit the number.
57  *
58  * Returns 0 on success, else EINVAL.
59  */
60 static int AppendNumberToString(char* str, int num) {
61     if (str == NULL || num >= 1000) {
62         return EINVAL;
63     }
64
65     char num_str[4];
66     memset(num_str, 0, 4);
67     if (num < 10) {
68         num_str[0] = num + '0';
69     }
70     else if (num < 100) {
71         num_str[0] = (num / 10) + '0';
72         num_str[1] = (num % 10) + '0';
73     }
74     else {
75         num_str[0] = (num / 100) + '0';
76         num_str[1] = ((num / 10) % 10) + '0';
77         num_str[2] = (num % 10) + '0';
78     }
79
80     strcat(str, num_str);
81     return 0;
82 }
83
84 /**
85  * Returns the name the next-mounted filesystem should receive (e.g. drv0,
86  * drv1, etc.). Each call to this function will return a different string. The
87  * caller is responsible for freeing the returned string.
88  *
89  * @return A caller-free string representing the drive name.
90  *
91  * @maxirql IRQL_SCHEDULER
92  */
93 char* GenerateNewMountedDiskName() {
94     MAX_IRQL(IRQL_SCHEDULER);
95
96     char name[16];
97     strcpy(name, "drv");
98
99     AcquireSpinlockIrql(&type_table_lock);
100     int disk_num = next_mounted_disk_num++;
101     ReleaseSpinlockIrql(&type_table_lock);
102
103     AppendNumberToString(name, disk_num);
104     return strdup(name);
105 }
106
107 /**
108  * Returns the name the next-mounted raw disk should receive, based on its type (e.g. raw-hd0, raw-hd1,
109  * raw-fd0). Each call to this function will return a different string. The caller is responsible for
110  * freeing the returned string.
111  *
112  * @param type The type of disk (one of DISKUTIL_TYPE ...)
113  * @return The caller-free string representing the drive name.
114  *
115  * @maxirql IRQL_SCHEDULER
116  */
117 char* GenerateNewRawDiskName(int type) {
118     MAX_IRQL(IRQL_SCHEDULER);
119
120     char name[16];
121     strcpy(name, "raw-");
122
123     if (type >= _DISKUTIL_NUM_TYPES || type < 0) {
124         type = DISKUTIL_TYPE_OTHER;
125     }
126
127     strcat(name, type_strings[type]);
128
129     AcquireSpinlockIrql(&type_table_lock);
130     int disk_num = type_table[type]++;
131     ReleaseSpinlockIrql(&type_table_lock);
132
133     AppendNumberToString(name, disk_num);
134     return strdup(name);
135 }
136
137 /**
138  * Generates and returns the name of a partition from its partition index within a drive (e.g. part0, part1)
139  * The caller is responsible for freeing the returned string.
140  */
141 static char* GetPartitionNameString(int index) {
142     char name[16];
143     strcpy(name, "part");
144     AppendNumberToString(name, index);
145     return strdup(name);
146 }
147
148 /**
149  * Given a disk, this function detects, creates and mounts partitions on that disk.
150  * For each detected partition, the filesystem is also detected, and that is mounted if it exists.
151  * If the disk has no partitions, a 'whole disk partition' will be created, the filesystem will
152  * still be detected. This should only be called once per disk, after it has been initialised.

```

```

153 *
154 * @param disk The disk to scan for partitions
155 *
156 * @maxirq1 IRQ1_STANDARD
157 */
158 void CreateDiskPartitions(struct open_file* disk) {
159     EXACT_IRQ1(IRQ1_STANDARD);
160
161     struct open_file** partitions = GetPartitionsForDisk(disk);
162
163     if (partitions == NULL || partitions[0] == NULL) {
164         struct stat st = disk->node->stat;
165         struct vnode* whole_disk_partition = CreatePartition(disk, 0, st.st_size, 0, st.st_blksize, 0, false)->node;
166         VnodeOpCreate(disk->node, whole_disk_partition, GetPartitionNameString(0), 0, 0);
167         return;
168     }
169
170     for (int i = 0; partitions[i]; ++i) {
171         struct vnode* partition = partitions[i]->node;
172         char* str = GetPartitionNameString(i);
173         VnodeOpCreate(disk->node, partition, str, 0, 0);
174     }
175
176
177 void InitDiskPartitionHelper(struct disk_partition_helper* helper) {
178     helper->num_partitions = 0;
179 }
180
181 int DiskFollowHelper(struct disk_partition_helper* helper, struct vnode** out, const char* name) {
182     for (int i = 0; i < helper->num_partitions; ++i) {
183         if (!strcmp(helper->partition_names[i], name)) {
184             *out = helper->partitions[i];
185             return 0;
186         }
187     }
188
189     *out = NULL;
190     return EINVAL;
191 }
192
193 int DiskCreateHelper(struct disk_partition_helper* helper, struct vnode** in, const char* name) {
194     if (helper->num_partitions == MAX_PARTITIONS_PER_DISK) {
195         return EINVAL;
196     }
197
198     helper->partitions[helper->num_partitions] = *in;
199     helper->partition_names[helper->num_partitions] = (char*) name;
200     helper->num_partitions++;
201     return 0;
202 }

```

File: ./vfs/transfer.c

```

1
2 #include <transfer.h>
3 #include <assert.h>
4 #include <string.h>
5 #include <errno.h>
6 #include <virtual.h>
7 #include <arch.h>
8 #include <log.h>
9
10 static int ValidateCopy(const void* user_addr, size_t size, bool write) {
11     size_t initial_address = (size_t) user_addr;
12
13     /*
14      * Check if the memory range starts in user memory.
15      */
16     if (initial_address < ARCH_USER_AREA_BASE || initial_address >= ARCH_USER_AREA_LIMIT) {
17         return EINVAL;
18     }
19
20     size_t final_address = initial_address + size;
21
22     /*
23      * Check for overflow when the initial address and size are added. If it would overflow,
24      * we cancel the operation, as the user is obviously outside their range.
25      */
26     if (final_address < initial_address) {
27         return EINVAL;
28     }
29
30     /*
31      * Ensure the end of the memory range is in user memory. As user memory must be contiguous,
32      * this ensures the entire range is in user memory.
33      */
34     if (final_address < ARCH_USER_AREA_BASE || final_address >= ARCH_USER_AREA_LIMIT) {
35         return EINVAL;
36     }
37
38     /*
39      * We must now check if the USER (and possibly WRITE) bits are set on the memory pages
40      * being accessed.
41      */
42     size_t initial_page = initial_address / ARCH_PAGE_SIZE;
43     size_t pages = BytesToPages(size);
44
45     for (size_t i = 0; i < pages; ++i) {
46         size_t page = initial_page + i;
47         size_t permissions = GetVirtPermissions(page * ARCH_PAGE_SIZE);
48
49         if (permissions == 0) {
50             return EINVAL;
51         }
52
53         if (!(permissions & VM_READ)) {
54             return EINVAL;
55         }
56         if (!(permissions & VM_USER)) {
57             return EINVAL;
58         }
59         if (write && !(permissions & VM_WRITE)) {
60             return EINVAL;
61         }
62         if (write && (permissions & VM_EXEC)) {
63             return EINVAL;
64         }
65         if (write && (permissions & VM_EXEC)) {
66             return EINVAL;
67         }
68     }
69
70     return 0;
71 }
72
73 static int CopyIntoKernel(void* kernel_addr, const void* user_addr, size_t size) {
74     int status = ValidateCopy(user_addr, size, false);

```

```

75     if (status != 0) {
76         return status;
77     }
78     inline_memcpy(kernel_addr, user_addr, size);
79     return 0;
80 }
81
82
83 static int CopyOutOfKernel(const void* kernel_addr, void* user_addr, size_t size) {
84     int status = ValidateCopy(user_addr, size, true);
85     if (status != 0) {
86         return status;
87     }
88     inline_memcpy(user_addr, kernel_addr, size);
89     return 0;
90 }
91
92
93 int PerformTransfer(void* trusted_buffer, struct transfer* untrusted_buffer, uint64_t len) {
94     assert(trusted_buffer != NULL);
95     assert(untrusted_buffer != NULL && untrusted_buffer->address != NULL);
96     assert(untrusted_buffer->direction == TRANSFER_READ || untrusted_buffer->direction == TRANSFER_WRITE);
97
98     size_t amount_to_copy = MIN(len, untrusted_buffer->length_remaining);
99     if (amount_to_copy == 0) {
100         return 0;
101     }
102
103     if (untrusted_buffer->type == TRANSFER_INTRA_KERNEL) {
104         if (untrusted_buffer->direction == TRANSFER_READ) {
105             memmove(untrusted_buffer->address, trusted_buffer, amount_to_copy);
106         } else {
107             memmove(trusted_buffer, untrusted_buffer->address, amount_to_copy);
108         }
109     } else {
110         int result;
111         /*
112          * This is from the kernel's perspective of the operations.
113          */
114         if (untrusted_buffer->direction == TRANSFER_READ) {
115             result = CopyOutOfKernel((const void*) trusted_buffer, untrusted_buffer->address, amount_to_copy);
116         } else {
117             result = CopyIntoKernel(trusted_buffer, (const void*) untrusted_buffer->address, amount_to_copy);
118         }
119         if (result != 0) {
120             return result;
121         }
122     }
123
124     untrusted_buffer->length_remaining -= amount_to_copy;
125     untrusted_buffer->offset += amount_to_copy;
126     untrusted_buffer->address = ((uint8_t*) untrusted_buffer->address) + amount_to_copy;
127     return 0;
128 }
129
130
131 int WriteStringToUsermode(const char* trusted_string, char* untrusted_buffer, uint64_t max_length) {
132     struct transfer tr = CreateTransferWritingToUser(untrusted_buffer, max_length, 0);
133     int result;
134     /*
135      * Limit the size of the string by the maximum length. We use <, and a -1 in the other case,
136      * as we need to ensure the null terminator fits.
137      */
138     uint64_t size = strlen(trusted_string) < max_length ? strlen(trusted_string) : max_length - 1;
139     result = PerformTransfer((void*) trusted_string, &tr, size);
140     if (result != 0) {
141         return result;
142     }
143     uint8_t zero = 0;
144     return PerformTransfer(&zero, &tr, 1);
145 }
146
147
148 int ReadStringFromUsermode(char* trusted_buffer, const char* untrusted_string, uint64_t max_length) {
149     struct transfer tr = CreateTransferReadingFromUser(untrusted_string, max_length, 0);
150     size_t i = 0;
151     while (max_length-- > 1) {
152         char c;
153         int result = PerformTransfer(&c, &tr, 1);
154         if (result != 0) {
155             return result;
156         }
157         trusted_buffer[i++] = c;
158         if (c == 0) {
159             break;
160         }
161     }
162     trusted_buffer[i] = 0;
163     return 0;
164 }
165
166
167 int WriteWordToUsermode(size_t* location, size_t value) {
168     struct transfer io = CreateTransferWritingToUser(location, sizeof(size_t), 0);
169     int res = PerformTransfer(&value, &io, sizeof(size_t));
170     if (io.length_remaining != 0) {
171         return EINVAL;
172     }
173     return res;
174 }
175
176
177 int ReadWordFromUsermode(size_t* location, size_t* output) {
178     struct transfer io = CreateTransferReadingFromUser(location, sizeof(size_t), 0);
179     int res = PerformTransfer(output, &io, sizeof(size_t));
180     if (io.length_remaining != 0) {
181         return EINVAL;
182     }
183     return res;
184 }
185
186
187 static struct transfer CreateTransfer(void* addr, uint64_t length, uint64_t offset, int direction, int type) {
188     struct transfer trans;
189     trans.address = addr;
190     trans.direction = direction;
191     trans.length_remaining = length;
192     trans.offset = offset;
193     trans.type = type;
194     return trans;
195 }
196
197
198 struct transfer CreateKernelTransfer(void* addr, uint64_t length, uint64_t offset, int direction) {

```

```

205     return CreateTransfer(addr, length, offset, direction, TRANSFER_INTRA_KERNEL);
206 }
207
208 struct transfer CreateTransferWritingToUser(void* untrusted_addr, uint64_t length, uint64_t offset) {
209     /*
210      * When we "write to the user", we are doing so because the user is trying to "read" from the kernel.
211      * i.e. someone is doing an "untrusted read" of kernel data (i.e. a TRANSFER_READ).
212      */
213     return CreateTransfer(untrusted_addr, length, offset, TRANSFER_READ, TRANSFER_USERMODE);
214 }
215
216 struct transfer CreateTransferReadingFromUser(const void* untrusted_addr, uint64_t length, uint64_t offset) {
217     /*
218      * When we "read from the user", we are doing so because the user is trying to "write" to the kernel.
219      * i.e. as they are writing to the kernel, it is an "untrusted write" (i.e. a TRANSFER_WRITE).
220      */
221     return CreateTransfer((void*) untrusted_addr, length, offset, TRANSFER_WRITE, TRANSFER_USERMODE);
222 }

```

File: /vfs/filedes.c

```

1  #include <filedes.h>
2  #include <errno.h>
3  #include <string.h>
4  #include <common.h>
5  #include <semaphore.h>
6  #include <vfs.h>
7  #include <fcntl.h>
8  #include <log.h>
9  #include <heap.h>
10 #include <irq.h>
11
12 struct filedes_entry {
13     /*
14      * Set to NULL if this entry isn't in use.
15      */
16     struct open_file* file;
17
18     /*
19      * The only flag that can live here is FD_CLOEXEC. All other flags live on the filesystem
20      * level. This is because FD_CLOEXEC is a property of the file descriptor, not the underlying
21      * file itself. (This is important in how dup() works.).
22      */
23     * Note that we set FD_CLOEXEC == O_CLOEXEC.
24     */
25     int flags;
26 };
27
28 /*
29 * The table of all of the file descriptors in use by a process.
30 */
31 struct filedes_table {
32     struct semaphore* lock;
33     struct filedes_entry* entries;
34 };
35
36 struct filedes_table* CreateFileDescriptorTable(void) {
37     struct filedes_table* table = AllocHeap(sizeof(struct filedes_table));
38
39     table->lock = CreateMutex("filedes");
40     table->entries = AllocHeapEx(sizeof(struct filedes_entry) * MAX_FD_PER_PROCESS, HEAP_ALLOW_PAGING);
41
42     for (int i = 0; i < MAX_FD_PER_PROCESS; ++i) {
43         table->entries[i].file = NULL;
44     }
45
46     return table;
47 }
48
49 struct filedes_table* CopyFileDescriptorTable(struct filedes_table* original) {
50     struct filedes_table* new_table = CreateFileDescriptorTable();
51
52     AcquireMutex(original->lock, -1);
53     memcpy(new_table->entries, original->entries, sizeof(struct filedes_entry) * MAX_FD_PER_PROCESS);
54     ReleaseMutex(original->lock);
55
56     return new_table;
57 }
58
59 void DestroyFileDescriptorTable(struct filedes_table* table) {
60     AcquireMutex(table->lock, -1);
61
62     for (int i = 0; i < MAX_FD_PER_PROCESS; ++i) {
63         if (table->entries[i].file != NULL) {
64             CloseFile(table->entries[i].file);
65         }
66     }
67
68     ReleaseMutex(table->lock);
69     DestroyMutex(table->lock);
70 }
71
72 int CreateFileDescriptor(struct filedes_table* table, struct open_file* file, int* fd_out, int flags) {
73     if ((flags & O_CLOEXEC) != 0) {
74         return EINVAL;
75     }
76
77     AcquireMutex(table->lock, -1);
78
79     for (int i = 0; i < MAX_FD_PER_PROCESS; ++i) {
80         if (table->entries[i].file == NULL) {
81             table->entries[i].file = file;
82             table->entries[i].flags = flags;
83             ReleaseMutex(table->lock);
84             *fd_out = i;
85             return 0;
86         }
87     }
88
89     ReleaseMutex(table->lock);
90     return EMFILE;
91 }
92
93 int RemoveFileDescriptor(struct filedes_table* table, struct open_file* file) {
94     AcquireMutex(table->lock, -1);
95
96     for (int i = 0; i < MAX_FD_PER_PROCESS; ++i) {
97         if (table->entries[i].file == file) {
98             table->entries[i].file = NULL;
99             ReleaseMutex(table->lock);
100             return 0;
101         }
102     }
103
104     ReleaseMutex(table->lock);
105     return EINVAL;
106 }

```

```

107
108 int GetFileFromDescriptor(struct filedes_table* table, int fd, struct open_file** out) {
109     if (out == NULL || fd < 0 || fd >= MAX_FD_PER_PROCESS) {
110         *out = NULL;
111         return out == NULL ? EINVAL : EBADF;
112     }
113
114     AcquireMutex(table->lock, -1);
115     struct open_file* result = table->entries[fd].file;
116     ReleaseMutex(table->lock);
117
118     *out = result;
119     return result == NULL ? EBADF : 0;
120 }
121
122 int HandleFileDescriptorsOnExec(struct filedes_table* table) {
123     AcquireMutex(table->lock, -1);
124
125     for (int i = 0; i < MAX_FD_PER_PROCESS; ++i) {
126         if (table->entries[i].file != NULL) {
127             if (table->entries[i].flags & O_CLOEXEC) {
128                 struct open_file* file = table->entries[i].file;
129                 table->entries[i].file = NULL;
130                 int res = CloseFile(file);
131                 if (res != 0) {
132                     ReleaseMutex(table->lock);
133                     return res;
134                 }
135             }
136         }
137     }
138
139     ReleaseMutex(table->lock);
140     return 0;
141 }
142
143 int DuplicateFileDescriptor(struct filedes_table* table, int oldfd, int* newfd) {
144     AcquireMutex(table->lock, -1);
145
146     struct open_file* original_file;
147     int res = GetFileFromDescriptor(table, oldfd, &original_file);
148     if (res != 0 || original_file == NULL) {
149         ReleaseMutex(table->lock);
150         return EBADF;
151     }
152
153     for (int i = 0; i < MAX_FD_PER_PROCESS; ++i) {
154         if (table->entries[i].file == NULL) {
155             table->entries[i].file = original_file;
156             table->entries[i].flags = 0;
157             ReleaseMutex(table->lock);
158             *newfd = i;
159             return 0;
160         }
161     }
162
163     ReleaseMutex(table->lock);
164     return EMFILE;
165 }
166
167 int DuplicateFileDescriptor2(struct filedes_table* table, int oldfd, int newfd, int flags) {
168     if (flags & ~O_CLOEXEC) {
169         return EINVAL;
170     }
171
172     AcquireMutex(table->lock, -1);
173
174     struct open_file* original_file;
175     int res = GetFileFromDescriptor(table, oldfd, &original_file);
176
177     /*
178      * "If oldfd is not a valid file descriptor, then the call fails,
179      * and newfd is not closed."
180      */
181     if (res != 0 || original_file == NULL) {
182         ReleaseMutex(table->lock);
183         return EBADF;
184     }
185
186     /*
187      * "If oldfd is a valid file descriptor, and newfd has the same
188      * value as oldfd, then dup2() does nothing..."
189      */
190     if (oldfd == newfd) {
191         ReleaseMutex(table->lock);
192         return 0;
193     }
194
195     struct open_file* current_file;
196     res = GetFileFromDescriptor(table, oldfd, &current_file);
197     if (res == 0 && current_file != NULL) {
198         /*
199          * "If the file descriptor newfd was previously open, it is closed
200          * before being reused; the close is performed silently (i.e., any
201          * errors during the close are not reported by dup2())."
202          */
203         CloseFile(current_file);
204     }
205
206     table->entries[newfd].file = original_file;
207     table->entries[newfd].flags = flags;
208
209     ReleaseMutex(table->lock);
210     return 0;
211 }

```

File: ./vfs/vnode.c

```

1  #include <vnode.h>
2  #include <spinlock.h>
3  #include <assert.h>
4  #include <log.h>
5  #include <spinlock.h>
6  #include <errno.h>
7  #include <heap.h>
8  #include <dirent.h>
9  #include <irql.h>
10 #include <string.h>
11 #include <sys/stat.h>
12
13 /*
14  * Allocate and initialise a vnode. The reference count is initialised to 1.
15  */
16 struct vnode* CreateVnode(struct vnode_operations ops, struct stat st) {
17     struct vnode* node = AllocHeap(sizeof(struct vnode));
18     *node = (struct vnode) {

```

```

19     .ops = ops, .reference_count = 1, .data = NULL, .stat = st
20 };
21 InitSpinlock(&node->reference_count_lock, "vnode_refcnt", IRQL_SCHEDULER);
22 return node;
23 }
24
25 /*
26 * Cleanup and free an abstract file node.
27 */
28 static void DestroyVnode(struct vnode* node) {
29     /*
30     * The lock can't be held during this process, otherwise the lock will
31     * get freed before it is released (which is bad, as we must release it
32     * to get interrupts back on).
33     */
34     assert(node != NULL);
35     assert(node->reference_count == 0);
36     if (node->stat.st_nlink == 0) {
37         VnodeOpDelete(node);
38     }
39     FreeHeap(node);
40 }
41
42 /*
43 * Ensures that a vnode is valid.
44 */
45 static void CheckVnode(struct vnode* node) {
46     assert(node != NULL);
47
48     if (IsSpinlockHeld(&node->reference_count_lock)) {
49         assert(node->reference_count > 0);
50     } else {
51         AcquireSpinlockIrql(&node->reference_count_lock);
52         assert(node->reference_count > 0);
53         ReleaseSpinlockIrql(&node->reference_count_lock);
54     }
55 }
56
57 /*
58 * Increments a vnode's reference counter. Used whenever a vnode is 'given' to someone.
59 */
60 void ReferenceVnode(struct vnode* node) {
61     assert(node != NULL);
62
63     AcquireSpinlockIrql(&node->reference_count_lock);
64     node->reference_count++;
65     ReleaseSpinlockIrql(&node->reference_count_lock);
66 }
67
68 /*
69 * Decrements a vnode's reference counter, destorying it if it reaches zero.
70 * It should be called to free a vnode 'given' to use when it is no longer needed.
71 */
72 void DereferenceVnode(struct vnode* node) {
73     CheckVnode(node);
74     AcquireSpinlockIrql(&node->reference_count_lock);
75
76     assert(node->reference_count > 0);
77     node->reference_count--;
78
79     if (node->reference_count == 0) {
80         VnodeOpClose(node);
81
82         /*
83          * Must release the lock before we delete it so we can put interrupts back on
84          */
85         ReleaseSpinlockIrql(&node->reference_count_lock);
86
87         DestroyVnode(node);
88         return;
89     }
90
91     ReleaseSpinlockIrql(&node->reference_count_lock);
92 }
93
94
95 int VnodeOpCheckOpen(struct vnode* node, const char* name, int flags) {
96     CheckVnode(node);
97     if (node->ops.check_open == NULL) {
98         return 0;
99     }
100     return node->ops.check_open(node, name, flags);
101 }
102
103 int VnodeOpRead(struct vnode* node, struct transfer* io) {
104     CheckVnode(node);
105     if (node->ops.read == NULL || io->direction != TRANSFER_READ) {
106         return EINVAL;
107     }
108     return node->ops.read(node, io);
109 }
110
111 int VnodeOpWrite(struct vnode* node, struct transfer* io) {
112     CheckVnode(node);
113     if (node->ops.write == NULL || io->direction != TRANSFER_WRITE) {
114         return EINVAL;
115     }
116     return node->ops.write(node, io);
117 }
118
119 int VnodeOpIoctl(struct vnode* node, int command, void* buffer) {
120     CheckVnode(node);
121     if (node->ops.ioctl == NULL) {
122         return EINVAL;
123     }
124     return node->ops.ioctl(node, command, buffer);
125 }
126
127 int VnodeOpClose(struct vnode* node) {
128     /*
129     * Don't call CheckVnode, as that tests the reference count being non-zero,
130     * but it should be zero here.
131     */
132     if (node->reference_count != 0) {
133         return EINVAL;
134     }
135     if (node->ops.close == NULL) {
136         return 0;
137     }
138     return node->ops.close(node);
139 }
140
141 int VnodeOpCreate(struct vnode* node, struct vnode** out, const char* name, int flags, mode_t mode) {
142     CheckVnode(node);
143     if (node->ops.create == NULL) {
144         return EINVAL;
145     }
146     return node->ops.create(node, out, name, flags, mode);
147 }
148

```



```

149 int VnodeOpTruncate(struct vnode* node, off_t offset) {
150     CheckVnode(node);
151     if (node->ops.truncate == NULL) {
152         return EINVAL;
153     }
154     return node->ops.truncate(node, offset);
155 }
156
157 int VnodeOpFollow(struct vnode* node, struct vnode** new_node, const char* name) {
158     CheckVnode(node);
159     if (node->ops.follow == NULL) {
160         return ENOTDIR;
161     }
162     return node->ops.follow(node, new_node, name);
163 }
164
165 uint8_t VnodeOpDirentType(struct vnode* node) {
166     return IFTODT(node->stat.st_mode);
167 }
168
169 int VnodeOpWait(struct vnode* node, int flags, uint64_t timeout_ms) {
170     CheckVnode(node);
171     if (node->ops.wait == NULL) {
172         return 0;
173     }
174     return node->ops.wait(node, flags, timeout_ms);
175 }
176
177 int VnodeOpUnlink(struct vnode* node) {
178     CheckVnode(node);
179     if (node->ops.unlink == NULL) {
180         return EINVAL;
181     }
182     return node->ops.unlink(node);
183 }
184
185 int VnodeOpDelete(struct vnode* node) {
186     CheckVnode(node);
187     if (node->ops.delete == NULL) {
188         return EINVAL;
189     }
190     return node->ops.delete(node);
191 }

```

File: ./vfs/vfs.c

```

1  #include <vfs.h>
2  #include <spinlock.h>
3  #include <irq.h>
4  #include <log.h>
5  #include <assert.h>
6  #include <virtual.h>
7  #include <string.h>
8  #include <errno.h>
9  #include <dirent.h>
10 #include <heap.h>
11 #include <avl.h>
12 #include <fcntl.h>
13 #include <stackadt.h>
14
15 /*
16  * Try not to have non-static functions that return in any way a struct vnode*, as it
17  * probably means you need to use the reference/dereference functions.
18  */
19
20 /*
21  * Maximum length of a component of a filepath (e.g. an file/directory's individual name).
22  */
23 #define MAX_COMPONENT_LENGTH 128
24
25 /*
26  * Maximum total length of a path.
27  */
28 #define MAX_PATH_LENGTH 2000
29
30 /*
31  * Maximum number of symbolic links to dereference in a path before returning ELOOP.
32  */
33 #define MAX_LOOP 5
34
35
36 /*
37  * A structure for mounted devices and filesystems.
38  */
39 struct mounted_file {
40     /* The vnode representing the device / root directory of a filesystem */
41     struct open_file* node;
42
43     /* What the device / filesystem mount is called */
44     char* name;
45 };
46
47 static struct spinlock vfs_lock;
48 static struct avl_tree* mount_points = NULL;
49
50 int MountedDeviceComparator(void* a_, void* b_) {
51     struct mounted_file* a = a_;
52     struct mounted_file* b = b_;
53     return strcmp(a->name, b->name);
54 }
55
56 void InitVfs(void) {
57     InitSpinlock(&vfs_lock, "vfs", IRQL_SCHEDULER);
58     mount_points = AvlTreeCreate();
59     AvlTreeSetComparator(mount_points, MountedDeviceComparator);
60 }
61
62 static int CheckValidComponentName(const char* name)
63 {
64     assert(name != NULL);
65
66     if (name[0] == 0) {
67         return EINVAL;
68     }
69
70     for (int i = 0; name[i]; ++i) {
71         char c = name[i];
72
73         if (c == '/' || c == '\\\ ' || c == ':') {
74             return EINVAL;
75         }
76     }
77
78     return 0;
79 }
80
81 static int DoesMountPointExist(const char* name) {

```

```

82     assert(name != NULL);
83     assert(IsSpinlockHeld(&vfs_lock));
84
85     struct mounted_file target;
86     target.name = (char*) name;
87     if (AvlTreeContains(mount_points, &target)) {
88         return EEXIST;
89     }
90
91     return 0;
92 }
93
94 /*
95  * Given a filepath, and a pointer to an index within that filepath (representing where
96  * start searching), copies the next component into an output buffer of a given length.
97  * The index is updated to point to the start of the next component, ready for the next call.
98  *
99  * This also handles duplicated and trailing forward slashes.
100 */
101 static int GetPathComponent(const char* path, int* ptr, char* output_buffer, int max_output_length, char delimiter) {
102     int i = 0;
103
104     assert(path != NULL);
105     assert(ptr != NULL);
106     assert(max_output_length >= 1);
107     assert(delimiter != 0);
108     assert(output_buffer != NULL);
109     assert(strlen(path) >= 1);
110     assert(*ptr >= 0 && *ptr < (int) strlen(path));
111
112     /*
113      * These were meant to be caught at a higher level, so we can apply the current
114      * working directory or the current drive.
115      */
116     assert(path[0] != '/');
117     assert(path[0] != ':');
118
119     output_buffer[0] = 0;
120
121     while (path[*ptr] && path[*ptr] != delimiter) {
122         if (i >= max_output_length - 1) {
123             return ENAMETOOLONG;
124         }
125
126         /*
127          * Ensure we always have a null terminated string.
128          */
129         output_buffer[i++] = path[*ptr];
130         output_buffer[i] = 0;
131         (*ptr)++;
132     }
133
134     /*
135      * Skip past the delimiter (unless we are at the end of the string),
136      * as well as any trailing slashes (which could be after a slash delimiter, or
137      * after a colon).
138      */
139     if (path[*ptr]) {
140         do {
141             (*ptr)++;
142         } while (path[*ptr] == '/');
143     }
144
145     /*
146      * Ensure that there are no colons or backslashes in the filename itself.
147      */
148     return CheckValidComponentName(output_buffer);
149 }
150
151 static int GetFinalPathComponent(const char* path, char* output_buffer, int max_output_length) {
152     int path_ptr = 0;
153
154     int status = GetPathComponent(path, &path_ptr, output_buffer, max_output_length, ':');
155     if (status) {
156         return status;
157     }
158
159     while (path_ptr < (int) strlen(path)) {
160         status = GetPathComponent(path, &path_ptr, output_buffer, max_output_length, '/');
161         if (status) {
162             return status;
163         }
164     }
165
166     return 0;
167 }
168
169 /*
170  * Takes in a device name, without the colon, and returns its vnode.
171  * If no such device is mounted, it returns NULL.
172 */
173 static struct open_file* GetMountFromName(const char* name) {
174     assert(name != NULL);
175
176     if (mount_points == NULL) {
177         return NULL;
178     }
179
180     struct mounted_file target;
181     target.name = (char*) name;
182     struct mounted_file* mount = AvlTreeGet(mount_points, (void*) &target);
183     return mount->node;
184 }
185
186 int PAGEABLE_CODE_SECTION AddVfsMount(struct vnode* node, const char* name) {
187     MAX_IRQL(IRQL_PAGE_FAULT);
188
189     if (name == NULL || node == NULL) {
190         return EINVAL;
191     }
192
193     if (strlen(name) >= MAX_COMPONENT_LENGTH) {
194         return ENAMETOOLONG;
195     }
196
197     int status = CheckValidComponentName(name);
198     if (status != 0) {
199         return status;
200     }
201
202     AcquireSpinlockIrql(&vfs_lock);
203
204     if (DoesMountPointExist(name) == EEXIST) {
205         ReleaseSpinlockIrql(&vfs_lock);
206         return EEXIST;
207     }
208
209     struct mounted_file* mount = AllocHeap(sizeof(struct mounted_file));
210     mount->name = strdup(name);
211     mount->node = CreateOpenFile(node, 0, 0, true, true);

```

```

212
213     AvlTreeInsert(mount_points, (void*) mount);
214
215     LogWriteSerial("MOUNTED TO THE VFS: %s\n", name);
216
217     ReleaseSpinlockIrql(&vfs_lock);
218     return 0;
219
220
221 int PAGEABLE_CODE_SECTION RemoveVfsMount(const char* name) {
222     MAX_IRQL(IRQL_PAGE_FAULT);
223
224     if (name == NULL) {
225         return EINVAL;
226     }
227
228     if (CheckValidComponentName(name) != 0) {
229         return EINVAL;
230     }
231
232     AcquireSpinlockIrql(&vfs_lock);
233
234     /*
235     * Scan through the mount table for the device
236     */
237     struct mounted_file target;
238     target.name = (char*) name;
239
240     struct mounted_file* actual = AvlTreeGet(mount_points, &target);
241     if (actual == NULL) {
242         ReleaseSpinlockIrql(&vfs_lock);
243         return ENODEV;
244     }
245
246     assert(!strcmp(actual->name, name));
247
248     /*
249     * Decrement the reference that was initially created way back in
250     * vfs_add_device in the call to dev_create_vnode (the vnode dereference),
251     * and then the open file that was created alongside it.
252     */
253     DereferenceVnode(actual->node->node);
254     DereferenceOpenFile(actual->node);
255
256     AvlTreeDelete(mount_points, actual);
257     FreeHeap(actual->name);
258
259     ReleaseSpinlockIrql(&vfs_lock);
260     return 0;
261 }
262
263 static void CleanupVnodeStack(struct stack_adt* stack) {
264     /*
265     * We need to call dereference on each vnode in the stack before we
266     * can call StackAdtDestroy.
267     */
268     while (StackAdtSize(stack) > 0) {
269         struct vnode* node = StackAdtPop(stack);
270         DereferenceVnode(node);
271     }
272
273     StackAdtDestroy(stack);
274 }
275
276 /*
277 * Given an absolute filepath, returns the vnode representing
278 * the file, directory or device.
279 *
280 * Should be used carefully, as the reference count is incremented.
281 */
282 static int GetVnodeFromPath(const char* path, struct vnode** out, bool want_parent) {
283     assert(path != NULL);
284     assert(out != NULL);
285
286     if (strlen(path) == 0) {
287         return EINVAL;
288     }
289     if (strlen(path) >= MAX_PATH_LENGTH) {
290         return ENAMETOOLONG;
291     }
292
293     int path_ptr = 0;
294     char component_buffer[MAX_COMPONENT_LENGTH];
295
296     int err = GetPathComponent(path, &path_ptr, component_buffer, MAX_COMPONENT_LENGTH, '/');
297     if (err != 0) {
298         return err;
299     }
300
301     struct open_file* current_file = GetMountFromName(component_buffer);
302
303     /*
304     * No root device found, so we can't continue.
305     */
306     if (current_file == NULL || current_file->node == NULL) {
307         return ENODEV;
308     }
309
310     struct vnode* current_vnode = current_file->node;
311
312     /*
313     * This will be dereferenced either as we go through the loop, or
314     * after a call to vfs_close (this function should only be called
315     * by vfs_open).
316     */
317     ReferenceVnode(current_vnode);
318
319     char component[MAX_COMPONENT_LENGTH + 1];
320
321     /*
322     * To go back to a parent directory, we need to keep track of the previous component.
323     * As we can go back through many parents, we must keep track of all of them, hence we
324     * use a stack to store each vnode we encounter. We will not dereference the vnodes
325     * on the stack until the end using cleanup_vnode_stack.
326     */
327     struct stack_adt* previous_components = StackAdtCreate();
328
329     /*
330     * Iterate over the rest of the path.
331     */
332     while (path_ptr < (int) strlen(path)) {
333         int status = GetPathComponent(path, &path_ptr, component, MAX_COMPONENT_LENGTH, '/');
334         if (status != 0) {
335             DereferenceVnode(current_vnode);
336             CleanupVnodeStack(previous_components);
337             return status;
338         }
339
340         if (!strcmp(component, ".") || !strcmp(component, "..")) {
341             /*

```

```

342  * This doesn't change where we point to.
343  */
344  continue;
345
346  } else if (!strcmp(component, "..")) {
347      if (StackAdtSize(previous_components) == 0) {
348          /*
349           * We have reached the root. Going 'further back' than the root
350           * on Linux just keeps us at the root, so don't do anything here.
351           */
352      } else {
353          /*
354           * Pop the previous component and use it.
355           */
356          current_vnode = StackAdtPop(previous_components);
357      }
358      continue;
359  }
360
361  /*
362   * Use a separate pointer so that both inputs don't point to the same
363   * location. vnode follow either increments the reference count or creates
364   * a new vnode with a count of one.
365   */
366  struct vnode* next_vnode = NULL;
367  status = VnodeOpFollow(current_vnode, &next_vnode, component);
368  if (status != 0) {
369      DereferenceVnode(current_vnode);
370      CleanupVnodeStack(previous_components);
371      return status;
372  }
373
374  /*
375   * We have a component that can be backtracked to, so add it to the stack.
376   *
377   * Also note that vnode follow adds a reference count, so current vnode
378   * needs to be dereferenced. Conveniently, all components that need to be
379   * put on the stack also need dereferencing, and vice versa.
380   *
381   * The final vnode we find will not be added to the stack and dereferenced
382   * as we won't get here.
383   */
384  StackAdtPush(previous_components, current_vnode);
385  current_vnode = next_vnode;
386
387  int status = 0;
388
389  if (want_parent) {
390      /*
391       * Operations that require us to get the parent don't work if we are already
392       * at the root.
393       */
394      if (StackAdtSize(previous_components) == 0) {
395          status = EINVAL;
396      } else {
397          *out = StackAdtPop(previous_components);
398      }
399      CleanupVnodeStack(previous_components);
400      return status;
401  }
402
403  int RemoveFileOrDirectory(const char* path, bool rmdir) {
404      struct vnode* node;
405      int res = GetVnodeFromPath(path, &node, false);
406      if (res != 0) {
407          return res;
408      }
409      bool is_dir = IFTODT(node->stat.st_mode) == DT_DIR;
410      if (rmdir && !is_dir) return ENOTDIR;
411      if (!rmdir && is_dir) return EISDIR;
412
413      if (rmdir) {
414          res = VnodeOpDelete(node);
415      } else {
416          res = node->stat.st_nlink > 0 ? VnodeOpUnlink(node) : ENOENT;
417      }
418      DereferenceVnode(node);
419      return res;
420  }
421
422  int OpenFile(const char* path, int flags, mode_t mode, struct open_file* out) {
423      EXACT_IQQL(IQQL_STANDARD);
424      if (path == NULL || out == NULL || strlen(path) <= 0) {
425          return EINVAL;
426      }
427
428      int status;
429      char name[MAX_COMPONENT_LENGTH + 1];
430      status = GetFinalPathComponent(path, name, MAX_COMPONENT_LENGTH);
431      if (status) {
432          return status;
433      }
434
435      /*
436       * Grab the vnode from the path.
437       */
438      struct vnode* node;
439
440      /*
441       * Lookup a (hopefully) existing file.
442       */
443      status = GetVnodeFromPath(path, &node, false);
444
445      if (flags & O_CREAT) {
446          if (status == ENOENT) {
447              /*
448               * Get the parent folder.
449               */
450              status = GetVnodeFromPath(path, &node, true);
451              if (status) {
452                  return status;
453              }
454
455              struct vnode* child;
456              status = VnodeOpCreate(node, &child, name, flags, mode);
457              DereferenceVnode(node);
458
459              if (status) {

```

```

472     return status;
473 }
474
475 node = child;
476
477 } else if (flags & O_EXCL) {
478     /*
479     * The file already exists (as we didn't get EEXIST), but we were passed O_EXCL so we
480     * must give an error. If O_EXCL isn't passed, then O_CREAT will just open the existing file.
481     */
482     return EEXIST;
483 }
484
485 } else if (status != 0) {
486     return status;
487 }
488
489 status = VnodeOpCheckOpen(node, name, flags & (O_ACCMODE | O_NONBLOCK));
490 if (status) {
491     DereferenceVnode(node);
492     return status;
493 }
494
495 bool can_read = (flags & O_ACCMODE) != O_WRONLY;
496 bool can_write = (flags & O_ACCMODE) != O_RDONLY;
497 uint8_t dirent_type = VnodeOpDirentType(node);
498
499 if (dirent_type == DT_DIR && can_write) {
500     /*
501     * You cannot write to a directory. This also prevents truncation.
502     */
503     DereferenceVnode(node);
504     return EISDIR;
505 }
506
507 if ((flags & O_TRUNC) && dirent_type == DT_REG) {
508     if (!can_write)
509         status = VnodeOpTruncate(node, 0);
510     if (status)
511         return status;
512 }
513 return ENOSYS;
514 } else {
515     return EINVAL;
516 }
517 }
518
519 /* TODO: clear out the flags that don't normally get saved */
520
521 // TODO: may need to actually have a VnodeOpOpen, for things like FatFS.
522
523 *out = CreateOpenFile(node, mode, flags, can_read, can_write);
524 return 0;
525
526
527
528 static int FileAccess(struct open_file* file, struct transfer* io, bool write) {
529     EXACT_IQRL(IQRL_STANDARD);
530
531     if (io == NULL || io->address == NULL || file == NULL || file->node == NULL) {
532         return EINVAL;
533     }
534     if ((!write && !file->can_read) || (write && !file->can_write)) {
535         return EBADF;
536     }
537
538     if (file->flags & O_NONBLOCK) {
539         int block_status = VnodeOpWait(file->node, (write ? VNODE_WAIT_WRITE : VNODE_WAIT_READ) | VNODE_WAIT_NON_BLOCK, 0);
540         if (block_status != 0) {
541             return block_status;
542         }
543     }
544
545     if (write) {
546         return VnodeOpWrite(file->node, io);
547     } else {
548         return VnodeOpRead(file->node, io);
549     }
550 }
551
552 int ReadFile(struct open_file* file, struct transfer* io) {
553     return FileAccess(file, io, false);
554 }
555
556 int WriteFile(struct open_file* file, struct transfer* io) {
557     return FileAccess(file, io, true);
558 }
559
560 int CloseFile(struct open_file* file) {
561     EXACT_IQRL(IQRL_STANDARD);
562
563     if (file == NULL || file->node == NULL) {
564         return EINVAL;
565     }
566
567     DereferenceVnode(file->node);
568     DereferenceOpenFile(file);
569     return 0;
570 }

```

File: ./DS_Store

[binary]

File: ./util/Log.c

```

1
2 #include <common.h>
3 #include <log.h>
4
5 #define REAL_HW 0
6
7 __attribute__((no_instrument_function)) static void IntToStr(uint32_t i, char* output, int base)
8 {
9     const char* digits = "0123456789ABCDEF";
10
11     /*
12     * Work out where the end of the string is (this is based on the number).
13     * Using the do...while ensures that we always get at least one digit
14     * (i.e. ensures a 0 is printed if the input was 0).
15     */
16     uint32_t shifter = i;

```

```

17 do {
18     ++output;
19     shifter /= base;
20 } while (shifter);
21
22 /* Put in the null terminator. */
23 *output = '\0';
24
25 /*
26  * Now fill in the digits back-to-front.
27  */
28 do {
29     *--output = digits[i % base];
30     i /= base;
31 } while (i);
32 }
33
34 #if REAL_HW == 0
35 __attribute__((no_instrument_function)) static void outb(uint16_t port, uint8_t value)
36 {
37     asm volatile ("outb %0, %1" : : "a"(value), "Nd"(port));
38 }
39
40 __attribute__((no_instrument_function)) static uint8_t inb(uint16_t port)
41 {
42     uint8_t value;
43     asm volatile ("inb %1, %0"
44     : "=a"(value)
45     : "Nd"(port));
46     return value;
47 }
48 #endif
49
50 __attribute__((no_instrument_function)) static void logcnv(char c, bool screen)
51 {
52     if (screen) {
53         DbgScreenPutchar(c);
54     }
55     #if REAL_HW == 0
56     while ((inb(0x3F8 + 5) & 0x20) == 0) {
57         ;
58     }
59     outb(0x3F8, c);
60 #endif
61 }
62
63 __attribute__((no_instrument_function)) static void logsnv(char* a, bool screen)
64 {
65     while (*a) logcnv(*a++, screen);
66 }
67
68 __attribute__((no_instrument_function)) static void log_intnv(uint32_t i, int base, bool screen)
69 {
70     char str[12];
71     IntToStr(i, str, base);
72     logsnv(str, screen);
73 }
74
75 __attribute__((no_instrument_function)) static void LogWriteSerialVa(const char* format, va_list list, bool screen) {
76     if (format == NULL) {
77         format = "NULL";
78     }
79
80     int i = 0;
81
82     while (format[i]) {
83         if (format[i] == '%') {
84             switch (format[++i]) {
85                 case '%':
86                     logcnv('%', screen); break;
87                 case 'c':
88                     logcnv(va_arg(list, int), screen); break;
89                 case 's':
90                     logsnv(va_arg(list, char*), screen); break;
91                 case 'd':
92                     log_intnv(va_arg(list, signed), 10, screen); break;
93                 case 'x':
94                 case 'X':
95                     log_intnv(va_arg(list, unsigned), 16, screen); break;
96                 case 'l':
97                 case 'L':
98                     log_intnv(va_arg(list, unsigned long long), 16, screen); break;
99                 case 'u':
100                     log_intnv(va_arg(list, unsigned), 10, screen); break;
101                 default:
102                     logcnv('%', screen);
103                     logcnv(format[i], screen);
104                     break;
105             }
106         } else {
107             logcnv(format[i], screen);
108         }
109         i++;
110     }
111 }
112
113 __attribute__((no_instrument_function)) void LogWriteSerial(const char* format, ...)
114 {
115     va_list list;
116     va_start(list, format);
117     LogWriteSerialVa(format, list, false);
118     va_end(list);
119 }
120
121 __attribute__((no_instrument_function)) void LogDeveloperWarning(const char* format, ...) {
122     va_list list;
123     va_start(list, format);
124     LogWriteSerial("\n!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!\n\n>>> KERNEL DEVELOPER WARNING:\n    ");
125     LogWriteSerialVa(format, list, false);
126     va_end(list);
127 }
128
129 __attribute__((no_instrument_function)) void DbgScreenPrintf(const char* format, ...) {
130     va_list list;
131     va_start(list, format);
132     LogWriteSerialVa(format, list, true);
133     va_end(list);
134 }

```

File: /util/unicode.c

```

1
2 #include <unicode.h>
3 #include <errno.h>
4
5 /*
6  * Converts UTF16 to UTF32. out_length is in/out - it must contain the output buffer's maximum length,

```

```

7  * but will receive the result's length on success. NOT NULL TERMINATED ON INPUT OR OUTPUT!!
8  */
9  int Utf16ToCodepoints(uint16_t* utf16, uint32_t* codepoints, int in_length, int* out_length) {
10     int in = 0;
11     int out = 0;
12
13     while (in < in_length) {
14         if (out == out_length) {
15             return ENAMETOOLONG;
16         }
17
18         uint32_t codepoint = utf16[in++];
19         if (codepoint >= 0xD800 && codepoint <= 0xDBFF) {
20             if (in == in_length) {
21                 return EINVAL;
22             }
23             uint16_t low_surrogate = utf16[in++];
24             if (low_surrogate >= 0xDC00 && low_surrogate <= 0xDFFF) {
25                 codepoint = (codepoint - 0xD800) * 0x400 + (low_surrogate - 0xDC00);
26             } else {
27                 return EINVAL;
28             }
29         } else if (codepoint >= 0xDC00 && codepoint <= 0xDFFF) {
30             return EINVAL;
31         }
32
33         if (codepoint >= 0xD800 && codepoint <= 0xDFFF) {
34             return EINVAL;
35         }
36     }
37     codepoints[out++] = codepoint;
38 }
39
40 *out_length = out;
41 return 0;
42 }
43
44 int Utf8ToCodepoints(uint8_t* utf8, uint32_t* codepoints, int in_length, int* out_length) {
45     int in = 0;
46     int out = 0;
47
48     while (in < in_length) {
49         if (out == out_length) {
50             return ENAMETOOLONG;
51         }
52
53         uint32_t codepoint = utf8[in++];
54         if ((codepoint >> 3) == 0x1E) {
55             if (in + 2 >= in_length) {
56                 return EINVAL;
57             }
58             codepoint &= 0x7;
59             codepoint <<= 6;
60             uint8_t next = utf8[in++];
61             if ((next >> 6) != 0x2) {
62                 return EINVAL;
63             }
64             codepoint |= next & 0x3F;
65             codepoint <<= 6;
66             next = utf8[in++];
67             if ((next >> 6) != 0x2) {
68                 return EINVAL;
69             }
70             codepoint |= next & 0x3F;
71             next = utf8[in++];
72             if ((next >> 6) != 0x2) {
73                 return EINVAL;
74             }
75             codepoint |= next & 0x3F;
76         } else if ((codepoint >> 4) == 0xE) {
77             if (in + 1 >= in_length) {
78                 return EINVAL;
79             }
80             codepoint &= 0xF;
81             codepoint <<= 6;
82             uint8_t next = utf8[in++];
83             if ((next >> 6) != 0x2) {
84                 return EINVAL;
85             }
86             codepoint |= next & 0x3F;
87             next = utf8[in++];
88             if ((next >> 6) != 0x2) {
89                 return EINVAL;
90             }
91             codepoint |= next & 0x3F;
92         } else if ((codepoint >> 5) == 0x6) {
93             if (in >= in_length) {
94                 return EINVAL;
95             }
96             codepoint &= 0x1F;
97             codepoint <<= 6;
98             uint8_t next = utf8[in++];
99             if ((next >> 6) != 0x2) {
100                 return EINVAL;
101             }
102             codepoint |= next & 0x3F;
103         } else if (codepoint >= 0x80) {
104             return EINVAL;
105         }
106
107         if (codepoint >= 0xD800 && codepoint <= 0xDFFF) {
108             return EINVAL;
109         }
110
111         codepoints[out++] = codepoint;
112     }
113
114 *out_length = out;
115 return 0;
116 }
117
118 int CodepointsToUtf16(uint32_t* codepoints, uint16_t* utf16, int in_length, int* out_length) {
119     int in = 0;
120     int out = 0;
121     while (in < in_length) {
122         if (out == out_length) {
123             return ENAMETOOLONG;
124         }
125
126         uint32_t codepoint = codepoints[in++];
127         if (codepoint >= 0xD800 && codepoint <= 0xDFFF) {
128             return EINVAL;
129         }
130     }
131 }

```

```

137     if (codepoint > 0x10FFFF) {
138         return EINVAL;
139     }
140
141     if (codepoint >= 0x10000) {
142         uint16_t high_surrogate = (codepoint / 0x400) + 0xD800;
143         uint16_t low_surrogate = (codepoint % 0x400) + 0xDC00;
144         utf16[out++] = high_surrogate;
145         if (out == *out_length) {
146             return ENAMETOOLONG;
147         }
148         utf16[out++] = low_surrogate;
149     } else {
150         utf16[out++] = codepoint;
151     }
152 }
153
154 *out_length = out;
155 return 0;
156 }
157
158 int CodepointsToUtf8(uint32_t* codepoints, uint8_t* utf8, int in_length, int* out_length) {
159     int in = 0;
160     int out = 0;
161     while (in < in_length) {
162         if (out == *out_length) {
163             return ENAMETOOLONG;
164         }
165
166         uint32_t codepoint = codepoints[in++];
167
168         if (codepoint >= 0xD800 && codepoint <= 0xDFFF) {
169             return EINVAL;
170         }
171
172         if (codepoint <= 0x7F) {
173             utf8[out++] = codepoint;
174         } else if (codepoint <= 0x7FF) {
175             if (out + 2 > *out_length) {
176                 return ENAMETOOLONG;
177             }
178             utf8[out++] = 0xC0 | (codepoint >> 6);
179             utf8[out++] = 0x80 | (codepoint & 0x3F);
180         } else if (codepoint <= 0xFFFF) {
181             if (out + 3 > *out_length) {
182                 return ENAMETOOLONG;
183             }
184             utf8[out++] = 0xE0 | (codepoint >> 12);
185             utf8[out++] = 0x80 | ((codepoint >> 6) & 0x3F);
186             utf8[out++] = 0x80 | (codepoint & 0x3F);
187         } else if (codepoint <= 0x10FFF) {
188             if (out + 4 > *out_length) {
189                 return ENAMETOOLONG;
190             }
191             utf8[out++] = 0xF0 | (codepoint >> 18);
192             utf8[out++] = 0x80 | ((codepoint >> 12) & 0x3F);
193             utf8[out++] = 0x80 | ((codepoint >> 6) & 0x3F);
194             utf8[out++] = 0x80 | (codepoint & 0x3F);
195         } else {
196             return EINVAL;
197         }
198     }
199
200     *out_length = out;
201     return 0;
202 }
203
204 *out_length = out;
205 return 0;
206 }
207
208

```

File: /util/video.c

```

1 #include <heap.h>
2 #include <assert.h>
3 #include <string.h>
4 #include <irq.h>
5 #include <video.h>
6 #include <virtual.h>
7 #include <errno.h>
8 #include <log.h>
9
10 static struct video_driver video_driver = {
11     .putchar = NULL,
12     .puts = NULL
13 };
14
15 void DeferPuts(void* v) {
16     video_driver.puts((char*) v);
17 }
18
19 void DeferPuchar(void* v) {
20     video_driver.putchar((char) (size_t) v);
21 }
22
23 void DbgScreenPuchar(char c) {
24     if (video_driver.putchar != NULL) {
25         DeferUntilIrql(IRQL_STANDARD, DeferPuchar, (void*) (size_t) c);
26     }
27 }
28
29 void DbgScreenPuts(char* s) {
30     if (video_driver.puts != NULL) {
31         DeferUntilIrql(IRQL_STANDARD, DeferPuts, s);
32     }
33 }
34
35 void InitVideoConsole(struct video_driver driver) {
36     video_driver = driver;
37 }

```

File: /util/panic.c


```

1
2 #include <panic.h>
3 #include <arch.h>
4 #include <log.h>
5 #include <errno.h>
6 #include <debug.h>
7 #include <irq.h>
8
9 static const char* message_table[_PANIC_HIGHEST_VALUE] = {
10 [PANIC_UNKNOWN] = "unknown",
11 [PANIC_IMPOSSIBLE_RETURN] = "impossible return",
12 [PANIC_MANUALLY_INITIATED] = "manually initiated",
13 [PANIC_UNIT_TEST_OK] = "unit test ok",
14 [PANIC_DRIVER_FAULT] = "driver fault",
15 [PANIC_OUT_OF_HEAP] = "out of heap",
16 [PANIC_OUT_OF_BOOTSTRAP_HEAP] = "out of bootstrap heap",
17 [PANIC_HEAP_REQUEST_TOO_LARGE] = "heap request too large",
18 [PANIC_ASSERTION_FAILURE] = "assertion failure",
19 [PANIC_NO_MEMORY_MAP] = "no memory map",
20 [PANIC_NOT_IMPLEMENTED] = "not implemented",
21 [PANIC_INVALID_IRQ] = "invalid irq level",
22 [PANIC_SPINLOCK_WRONG_IRQ] = "spinlock wrong irq",
23 [PANIC_PRIORITY_QUEUE] = "invalid operation on a priority queue",
24 [PANIC_OUT_OF_PHYS] = "no physical memory left",
25 [PANIC_LINKED_LIST] = "invalid operation on a linked list",
26 [PANIC_CANARY_DIED] = "kernel stack overflow detected",
27 [PANIC_SEMAPHORE_DESTROY_WHILE_HELD] = "tried to destroy a held semaphore",
28 [PANIC_SEM_BLOCK_WITHOUT_THREAD] = "semaphore acquisition would block before multithreading has started",
29 [PANIC_CANNOT_LOCK_MEMORY] = "cannot lock virtual memory",
30 [PANIC_THREAD_LIST] = "invalid operation on a thread list",
31 [PANIC_CANNOT_MALLOCS_WITHOUT_FAULTING] = "heap allocation cannot be completed without faulting",
32 [PANIC_NO_FILESYSTEM] = "no driver can access the boot filesystem",
33 [PANIC_BAD_KERNEL] = "kernel executable file is corrupted",
34 [PANIC_DISK_FAILURE_ON_SWAPFILE] = "failed to read from or write to swapfile",
35 [PANIC_NEGATIVE_SEMAPHORE] = "more semaphore releases than acquisitions",
36 [PANIC_NON_MASKABLE_INTERRUPT] = "unrecoverable hardware error",
37 [PANIC_UNHANDLED_KERNEL_EXCEPTION] = "unhandled cpu exception",
38 [PANIC_REQUIRED_DRIVER_MISSING_SYMBOL] = "driver is missing a required symbol",
39 [PANIC_REQUIRED_DRIVER_NOT_FOUND] = "a required driver could not be loaded",
40 [PANIC_NO_LOW_MEMORY] = "not enough conventional memory to satisfy request",
41 [PANIC_OUT_OF_SWAPFILE] = "out of swapfile space",
42 [PANIC_PROGRAM_LOADER] = "failed to load the program loader",
43 [PANIC_WAS_TRIED_TO_SELF_DESTRUCT] = "tried to destroy the current virtual address space",
44 [PANIC ACPI AML] = "acpi aml is invalid",
45 [PANIC_SPINLOCK_DOUBLE_ACQUISITION] = "spinlock attempted to be locked while currently held",
46 [PANIC_SPINLOCK_RELEASED_BEFORE_ACQUIRED] = "spinlock attempted to be released while not held",
47 [PANIC_DOUBLE_FREE_DETECTED] = "heap memory freed twice",
48 [PANIC_CONFLICTING_ALLOCATION_REQUIREMENTS] = "conflicting heap allocation requirements"
49 };
50
51 static void (*graphical_panic_handler)(int, const char*) = NULL;
52
53 int SetGraphicalPanicHandler(void (*handler)(int, const char*)) {
54 if (graphical_panic_handler == NULL) {
55 graphical_panic_handler = handler;
56 return 0;
57 } else {
58 return EALREADY;
59 }
60 }
61
62
63 const char* GetPanicMessageFromCode(int code) {
64 return code < _PANIC_HIGHEST_VALUE ? message_table[code] : "";
65 }
66
67 [[noreturn]] void Panic(int code)
68 {
69 PanicEx(code, GetPanicMessageFromCode(code));
70 }
71
72 [[noreturn]] void PanicEx(int code, const char* message) {
73 LogWriteSerial("PANIC %d %s\n", code, message);
74 if (IsInTfwTest()) {
75 LogWriteSerial("in test.\n");
76 FinishedTfwTest(code);
77 ArchSetPowerState(ARCH_POWER_STATE_REBOOT);
78 }
79
80 RaiseIrql(IRQL_HIGH);
81 LogWriteSerial("\n\n *** KERNEL PANIC ***\n\n0x%X - %s\n", code, message);
82
83 if (graphical_panic_handler != NULL) {
84 graphical_panic_handler(code, message);
85 }
86
87 while (1) {
88 ArchDisableInterrupts();
89 ArchStallProcessor();
90 }
91 }

```

File: ./util/assert.c

```

1 #include <assert.h>
2 #include <panic.h>
3 #include <arch.h>
4 #include <log.h>
5 #include <debug.h>
6
7 _Noreturn void AssertionFail(const char* file, const char* line, const char* condition, const char* msg) {
8 ArchDisableInterrupts();
9 LogWriteSerial("Assertion failed: %s %s [%s: %s]\n", condition, msg, file, line);
10 Panic(PANIC_ASSERTION_FAILURE);
11 }

```

File: ./util/console.c

```

1
2 #include <console.h>
3 #include <pty.h>
4 #include <vfs.h>
5 #include <thread.h>
6 #include <virtual.h>
7 #include <string.h>
8 #include <log.h>
9
10 static struct vnode* console_master;
11 static struct vnode* console_subordinate;
12
13 static struct open_file* open_console_master;
14 static struct open_file* open_console_subordinate;
15
16 static bool console_initialised = false;
17
18 /*
19  * NOTE: the console only echos input when there's someone waiting for input
20  * (which should be fine most of the time - the only people input is when it's waiting for input!)
21  */
22
23 static void ConsoleDriverThread(void*) {
24     AddVfsMount(console_subordinate, "con");
25
26     while (true) {
27         char c;
28         struct transfer tr = CreateKernelTransfer(&c, 1, 0, TRANSFER_READ);
29         ReadFile(open_console_master, &tr);
30         DbgScreenPutchar(c);
31     }
32 }
33
34 void InitConsole(void) {
35     CreatePseudoTerminal(&console_master, &console_subordinate);
36     open_console_master = CreateOpenFile(console_master, 0, 0, true, true);
37     open_console_subordinate = CreateOpenFile(console_subordinate, 0, 0, true, true);
38     CreateThread(ConsoleDriverThread, NULL, GetVas(), "con");
39     console_initialised = true;
40 }
41
42 void SendKeystrokeConsole(char c) {
43     if (!console_initialised) return;
44
45     struct transfer tr = CreateKernelTransfer(&c, 1, 0, TRANSFER_WRITE);
46     WriteFile(open_console_master, &tr);
47 }
48
49 char GetcharConsole(void) {
50     if (!console_initialised) return 0;
51
52     char c;
53     struct transfer tr = CreateKernelTransfer(&c, 1, 0, TRANSFER_READ);
54     ReadFile(open_console_subordinate, &tr);
55     return c;
56 }
57
58 void PutcharConsole(char c) {
59     if (!console_initialised) return;
60     struct transfer tr = CreateKernelTransfer(&c, 1, 0, TRANSFER_WRITE);
61     WriteFile(open_console_subordinate, &tr);
62 }
63
64 void PutsConsole(const char* s) {
65     for (int i = 0; s[i]; ++i) {
66         PutcharConsole(s[i]);
67     }
68 }

```

File: /util/driver.c

```

1 #include <arch.h>
2 #include <driver.h>
3 #include <semaphore.h>
4 #include <irq.h>
5 #include <avl.h>
6 #include <string.h>
7 #include <vfs.h>
8 #include <errno.h>
9 #include <heap.h>
10 #include <panic.h>
11 #include <log.h>
12 #include <stdlib.h>
13 #include <virtual.h>
14 #include <fcntl.h>
15 #include <ctype.h>
16 #include <assert.h>
17
18 static struct semaphore* driver_table_lock;
19 static struct semaphore* symbol_table_lock;
20 static struct avl_tree* loaded_drivers;
21 static struct avl_tree* symbol_table;
22
23 struct symbol {
24     const char* name;
25     size_t addr;
26 };
27
28 struct loaded_driver {
29     char* filename;
30     size_t relocation_point;
31     struct quick_relocation_table* quick_relocation_table;
32 };
33
34 static int SymbolComparator(void* a_, void* b_) {
35     struct symbol* a = a_;
36     struct symbol* b = b_;
37     return strcmp(a->name, b->name);
38 }
39
40 static int DriverTableComparatorByRelocationPoint(void* a_, void* b_) {
41     struct loaded_driver* a = a_;
42     struct loaded_driver* b = b_;
43     return COMPARE_SIGN(a->relocation_point, b->relocation_point);
44 }
45
46 static int DriverTableComparatorByName(void* a_, void* b_) {
47     struct loaded_driver* a = a_;
48     struct loaded_driver* b = b_;
49     return strcmp(a->filename, b->filename);
50 }
51
52 static int QuickRelocationTableComparator(const void* a_, const void* b_) {
53     struct quick_relocation a = *((struct quick_relocation*) a_);
54     struct quick_relocation b = *((struct quick_relocation*) b_);
55     return COMPARE_SIGN(a.address, b.address);
56 }

```

```

57
58 static size_t GetDriverAddressWithLockHeld(const char* name) {
59     struct loaded_driver dummy = {.filename = (char*) name};
60
61     AvlTreeSetComparator(loaded_drivers, DriverTableComparatorByName);
62     struct loaded_driver* res = AvlTreeGet(loaded_drivers, &dummy);
63     if (res == NULL) {
64         return 0;
65     }
66     return res->relocation_point;
67 }
68
69 static struct loaded_driver* GetDriverFromAddress(size_t relocation_point) {
70     AcquireMutex(driver_table_lock, -1);
71
72     AvlTreeSetComparator(loaded_drivers, DriverTableComparatorByRelocationPoint);
73     struct loaded_driver dummy = {.relocation_point = relocation_point};
74     struct loaded_driver* res = AvlTreeGet(loaded_drivers, &dummy);
75
76     ReleaseMutex(driver_table_lock);
77     return res;
78 }
79
80 size_t GetDriverAddress(const char* name) {
81     EXACT_IRQL(IRQL_STANDARD);
82
83     AcquireMutex(driver_table_lock, -1);
84     size_t res = GetDriverAddressWithLockHeld(name);
85     ReleaseMutex(driver_table_lock);
86     return res;
87 }
88
89 void InitSymbolTable(void) {
90     driver_table_lock = CreateMutex("drv table");
91     symbol_table_lock = CreateMutex("sym table");
92
93     loaded_drivers = AvlTreeCreate();
94     symbol_table = AvlTreeCreate();
95     AvlTreeSetComparator(symbol_table, SymbolComparator);
96
97     struct open_file* kernel_file;
98     if (!OpenFile("sys:/kernel.exe", O_RDONLY, 0, &kernel_file)) {
99         Panic(PANIC_NO_FILESYSTEM);
100     }
101     ArchLoadSymbols(kernel_file, 0);
102     CloseFile(kernel_file);
103 }
104
105 static bool DoesSymbolContainIllegalCharacters(const char* symbol) {
106     for (int i = 0; symbol[i]; ++i) {
107         if (!isalnum(symbol[i]) && symbol[i] != '_' ) {
108             return true;
109         }
110     }
111     return strlen(symbol) == 0;
112 }
113
114 void AddSymbol(const char* symbol, size_t address) {
115     EXACT_IRQL(IRQL_STANDARD);
116
117     if (DoesSymbolContainIllegalCharacters(symbol)) {
118         return;
119     }
120
121     struct symbol* entry = AllocHeap(sizeof(struct symbol));
122     entry->name = strdup(symbol);
123     entry->addr = address;
124
125     AcquireMutex(symbol_table_lock, -1);
126     if (AvlTreeContains(symbol_table, entry)) {
127         /*
128          * The kernel has some symbols declared 'static' to file scope, with
129          * duplicate names (e.g. in /dev each file has its own 'Stat'). These
130          * get exported for some reason so we end up with duplicate names. We
131          * must ignore these to avoid AVL issues. They are safe to ignore, as
132          * they were meant to be 'static' anyway.
133          */
134         FreeHeap(entry);
135     } else {
136         AvlTreeInsert(symbol_table, entry);
137     }
138     ReleaseMutex(symbol_table_lock);
139 }
140
141 /*
142 * Do not name a (global) function pointer you receive from this the same as the
143 * actual function - this may cause issues with duplicate symbols.
144 *
145 * e.g. don't do this at global scope:
146 * void (*MyFunc)(void) = GetSymbolAddress("MyFunc");
147 *
148 * Instead, do void (*_MyFunc)(void) or something.
149 */
150 size_t GetSymbolAddress(const char* symbol) {
151     EXACT_IRQL(IRQL_STANDARD);
152
153     struct symbol dummy = {.name = symbol};
154
155     AcquireMutex(symbol_table_lock, -1);
156     struct symbol* result = AvlTreeGet(symbol_table, &dummy);
157     ReleaseMutex(symbol_table_lock);
158
159     if (result == NULL) {
160         return 0;
161     } else {
162         assert(!strcmp(result->name, symbol));
163         return result->addr;
164     }
165 }
166
167 static int LoadDriver(const char* name) {
168     struct open_file* file;
169     int res;
170     if ((res = OpenFile(name, O_RDONLY, 0, &file))) {
171         return res;
172     }
173
174     struct loaded_driver* drv = AllocHeap(sizeof(struct loaded_driver));
175     drv->filename = strdup_pageable(name);
176     drv->quick_relocation_table = NULL;
177     if ((res = ArchLoadDriver(&drv->relocation_point, file, &drv->quick_relocation_table))) {
178         return res;
179     }
180 }

```

```

187     assert(drv->quick_relocation_table != NULL);
188
189     AvlTreeInsert(loaded_drivers, drv);
190     ArchLoadSymbols(file, drv->relocation_point - 0xD0000000); // TODO: ??? GET RID OF ARCH SPECIFIC DETAILS (0xD0000000)
191     return 0;
192 }
193
194 int RequireDriver(const char* name) {
195     EXACT_IRQL(IRQL_STANDARD);
196
197     LogWriteSerial("Requiring driver: %s\n", name);
198
199     AcquireMutex(driver_table_lock, -1);
200
201     if (GetDriverAddressWithLockHeld(name) != 0) {
202         ReleaseMutex(driver_table_lock);
203
204         /*
205          * Not an error that it's already loaded - ideally no one should care if it has already been loaded
206          * or not. Hence we give 0 (success), not EALREADY.
207          */
208         return 0;
209     }
210
211     int res = LoadDriver(name);
212     ReleaseMutex(driver_table_lock);
213     return res;
214 }
215
216 void SortQuickRelocationTable(struct quick_relocation_table* table) {
217     struct quick_relocation* entries = table->entries;
218     qsort_pageable((void*) entries, table->used_entries, sizeof(struct quick_relocation), QuickRelocationTableComparator);
219 }
220
221 void AddToQuickRelocationTable(struct quick_relocation_table* table, size_t addr, size_t val) {
222     assert(table->used_entries < table->total_entries);
223     table->entries[table->used_entries].address = addr;
224     table->entries[table->used_entries].value = val;
225     table->used_entries++;
226 }
227
228 struct quick_relocation_table* CreateQuickRelocationTable(int count) {
229     struct quick_relocation_table* table = AllocHeap(sizeof(struct quick_relocation_table));
230     struct quick_relocation* entries = (struct quick_relocation*) MapVirt(0, 0, count * sizeof(struct quick_relocation), VM_READ | VM_WRITE, NULL, 0);
231     table->entries = entries;
232     table->used_entries = 0;
233     table->total_entries = count;
234     return table;
235 }
236
237 static int BinarySearchComparator(const void* a_, const void* b_) {
238     struct quick_relocation a = *((struct quick_relocation*) a_);
239     struct quick_relocation b = *((struct quick_relocation*) b_);
240
241     size_t page_a = a.address / ARCH_PAGE_SIZE;
242     size_t page_b = b.address / ARCH_PAGE_SIZE;
243
244     if (page_a > page_b) {
245         return 1;
246     } else if (page_a < page_b) {
247         return -1;
248     } else {
249         return 0;
250     }
251 }
252
253 static void ApplyRelocationsToPage(struct quick_relocation_table* table, size_t virtual) {
254     struct quick_relocation target;
255     target.address = virtual;
256     LogWriteSerial("ApplyRelocationsToPage: looking for 0x%X\n", virtual);
257
258     struct quick_relocation* entry = bsearch(&target, table->entries, table->used_entries, sizeof(struct quick_relocation), BinarySearchComparator);
259     if (entry == NULL) {
260         PanicEx(PANIC_ASSERTION_FAILURE, "quick relocation table doesn't contain lookup - bsearch or qsort is probably bugged");
261     }
262
263     /*
264      * As we only did a binary search to look for anything on that page, we might be halfway through the page!
265      * Move back through the entries until we find the start of page (or the first entry in the table).
266      */
267     while ((entry->address / ARCH_PAGE_SIZE == virtual / ARCH_PAGE_SIZE) ||
268            (entry->address - sizeof(size_t) + 1 == virtual / ARCH_PAGE_SIZE) && entry != table->entries) {
269         entry --;
270     }
271
272     /*
273      * We went past the last one, so need to forward onto it again - unless we hit the start of the table.
274      */
275     if (entry != table->entries) {
276         entry ++;
277     }
278
279     /*
280      * We also need to lock these, as (and YES this has actually happened before):
281      * - if the relocation is on the boundary of the next one, and the next one is not present
282      * - and the next one is relocatable
283      * Then:
284      * - we make the next page writable
285      * - we try to do the straddle relocation
286      * - that causes a fault on the next page (not present)
287      * - that one is also relocatable, so it enables writing on that page
288      * - it finishes, and unmarks it as writable
289      * - we fail to do the relocation, as it is no longer writable
290      * By locking it first, we force the relocations on the second page to happen first, and then we can
291      * mark it as writable.
292      */
293     bool needs_write_low = (GetVirtPermissions(virtual) & VM_WRITE) == 0;
294     bool needs_write_high = false;
295     bool need_unlock_high = false;
296
297     if (needs_write_low) {
298         SetVirtPermissions(virtual, VM_WRITE, 0);
299     }
300
301     size_t final_address = table->entries[table->used_entries - 1].address;
302
303     while ((entry->address / ARCH_PAGE_SIZE == virtual / ARCH_PAGE_SIZE) ||
304            (entry->address - sizeof(size_t) + 1 == virtual / ARCH_PAGE_SIZE) {
305         if ((entry->address + sizeof(size_t) + 1) / ARCH_PAGE_SIZE != virtual / ARCH_PAGE_SIZE) {
306             need_unlock_high = !LockVirt(virtual + ARCH_PAGE_SIZE);
307             needs_write_high = (GetVirtPermissions(virtual + ARCH_PAGE_SIZE) & VM_WRITE) == 0;
308             if (needs_write_high) {
309                 SetVirtPermissions(virtual + ARCH_PAGE_SIZE, VM_WRITE, 0);
310             }
311         }
312     }

```

```

317     }
318
319     size_t* ref = (size_t*) entry->address;
320     *ref = entry->value;
321
322     if (entry->address == final_address) {
323         break;
324     }
325     entry += 1;
326 }
327
328 if (needs_write_low) {
329     SetVirtPermissions(virtual, 0, VM_WRITE);
330 }
331 if (needs_write_high) {
332     SetVirtPermissions(virtual + ARCH_PAGE_SIZE, 0, VM_WRITE);
333 }
334 if (need_unlock_high) {
335     UnlockVirt(virtual + ARCH_PAGE_SIZE);
336 }
337 }
338
339 void PerformDriverRelocationOnPage(struct vas*, size_t relocation_base, size_t virt) {
340     LogWriteSerial("PerformDriverRelocationOnPage A\n");
341     struct loaded_driver* drv = GetDriverFromAddress(relocation_base);
342     if (drv == NULL) {
343         PanicEx(PANIC_ASSERTION_FAILURE, "PerformDriverRelocationOnPage");
344     }
345     LogWriteSerial("PerformDriverRelocationOnPage B. driver at 0x%X\n", drv);
346
347     ApplyRelocationsToPage(drv->quick_relocation_table, virt);
348 }

```

File: ./Makefile

```

CC = /Users/alex/Desktop/NOS/toolchain/output/bin/i386-elf-gcc
AS = nasm
FAKE_CROSS_COMPILER = -m32 -l"." -include -Imachine/include -linclude/openlibm
COMPILE_FLAGS = -c -Os -fipa-icf -std=gnu2x -fno-sections -fno-strict-aliasing -DCOMPILER_KERNEL -Wall -Wextra -Wpedantic -Werror -Wcast-align=strict -Wpointer-arith -fmax-errors=5 -ffreestanding $(FAKE_CROSS_COMPILER) -Wno-infinite-recursion -fomit-frame-pointer
#-fno-finstrument-functions

LINK_FLAGS = -fuse-ld=gold -Wl,--icf=all -Wl,-Map=kernel.map -nostartfiles -nostdlib -lgcc

# Set by the higher level Makefile before calling us - changes depending on whether we are compiling the debug or release build
LINKER_STRIP =
CPPDEFINES =

COBJECTS = $(patsubst %.c, %.o, $(wildcard *.c) $(wildcard */*.c) $(wildcard */*/*.c) $(wildcard */*/*/*.c) $(wildcard */*/*/*/*.c) $(wildcard */*/*/*/*.c))
ASMOBJECTS = $(patsubst %.s, %.oo, $(wildcard *.s) $(wildcard */*.s) $(wildcard */*/*.s) $(wildcard */*/*/*.s) $(wildcard */*/*/*/*.s) $(wildcard */*/*/*/*.s))

oskernel: $(COBJECTS) $(ASMOBJECTS) $(HOBJECTS)
$(CC) -T machine/linker.ld -o KERNEL.EXE $^ $(LINK_FLAGS) $(LINKER_STRIP)
# rm -r include
# rm -r machine
rm Makefile

%.o: %.c
$(CC) $(CPPDEFINES) $(COMPILE_FLAGS) $^ -o $@
rm $^

%.oo: %.s
$(AS) -felf32 $^ -o $@
rm $^

```

File: ./include/filedes.h

```

#pragma once

#include <common.h>

#define MAX_FD_PER_PROCESS 1024

struct open_file;
struct filedes_table;

int CreateFileDescriptor(struct filedes_table* table, struct open_file* file, int* fd_out, int flags);
int RemoveFileDescriptor(struct filedes_table* table, struct open_file* file);
int GetFileFromDescriptor(struct filedes_table* table, int fd, struct open_file** out);

int HandleFileDescriptorsOnExec(struct filedes_table* table);

struct filedes_table* CreateFileDescriptorTable(void);
struct filedes_table* CopyFileDescriptorTable(struct filedes_table* original);
void DestroyFileDescriptorTable(struct filedes_table* table);

int DuplicateFileDescriptor(struct filedes_table* table, int oldfd, int* newfd);
int DuplicateFileDescriptor2(struct filedes_table* table, int oldfd, int newfd, int flags);

```

File: ./include/filesystem.h

```

#pragma once

#include <common.h>

struct open_file;

typedef int(*fs_mount_creator)(struct open_file*, struct open_file**);

void InitFilesystemTable(void);
int RegisterFilesystem(char* fs_name, fs_mount_creator mount);
int MountFilesystemForDisk(struct open_file* partition);

```

File: ./include/semaphore.h

```

#pragma once

```

```

#include <common.h>

struct semaphore;
struct thread;

#define SEM_BIG_NUMBER (1 << 30)

#define SEM_DONT_CARE 0
#define SEM_REQUIRE_ZERO 1
#define SEM_REQUIRE_FULL 2

struct semaphore* CreateSemaphore(const char* name, int max_count, int initial_count);
int AcquireSemaphore(struct semaphore* sem, int timeout_ms);
void ReleaseSemaphore(struct semaphore* sem);
int DestroySemaphore(struct semaphore* sem, int mode);
int GetSemaphoreCount(struct semaphore* sem);

#define CreateMutex(name) CreateSemaphore(name, 1, 0)
#define AcquireMutex(mtx, timeout_ms) AcquireSemaphore(mtx, timeout_ms)
#define ReleaseMutex(mtx) ReleaseSemaphore(mtx)
#define DestroyMutex(mtx) DestroySemaphore(mtx, SEM_REQUIRE_ZERO)

void CancelSemaphoreOfThread(struct thread* thr);

```

File: ./include/console.h

```

#pragma once

#include <common.h>

void InitConsole(void);
void SendKeystrokeConsole(char c);
char GetcharConsole(void);
void PutcharConsole(char c);
void PutsConsole(const char* s);

```

File: ./include/stackadt.h

```

#pragma once

#include <common.h>

struct stack_adt;

struct stack_adt* StackAdtCreate(void);
void StackAdtDestroy(struct stack_adt* stack);
void StackAdtPush(struct stack_adt* stack, void* data);
void* StackAdtPeek(struct stack_adt* stack);
void* StackAdtPop(struct stack_adt* stack);
int StackAdtSize(struct stack_adt* stack);

```

File: ./include/spinlock.h

```

#pragma once

#include <common.h>

struct thread;

struct spinlock {
    size_t lock;
    struct thread* owner;
    char name[16];
    int irq;
    int prev_irq;
};

void InitSpinlock(struct spinlock* lock, const char* name, int irq);

int AcquireSpinlockIrql(struct spinlock* lock);
void ReleaseSpinlockIrql(struct spinlock* lock);

void AcquireSpinlockDirect(struct spinlock* lock);
void ReleaseSpinlockDirect(struct spinlock* lock);

bool IsSpinlockHeld(struct spinlock* lock);

```

File: ./include/pty.h

```

#pragma once

#include <common.h>

struct vnode;

void CreatePseudoTerminal(struct vnode** master, struct vnode** subordinate);

```

File: ./include/debug.h

```

#include <debug/tfw.h>

```

File: ./include/progload.h

```

#pragma once

#include <common.h>

```

```
void InitProgramLoader(void);
int CopyProgramLoaderIntoAddressSpace(void);
```

File: ./include/irq.h

```
#pragma once

#include <assert.h>

struct irq_deferment {
    void (*handler)(void*);
    void* context;
};

/**
 * SIMPLE TABLE
 *
 * Can page fault? Can task switch? Can use drivers? Can have IRQs?
 * IRQL_STANDARD YES YES YES YES
 * IRQL_PAGE_FAULT SORT OF YES YES YES (only the page fault handler can generate a nested page fault, e.g. handling some COW stuff)
 * IRQL_SCHEDULER NO SORT OF YES YES (only the scheduler can make a task switch occur, others get postponed)
 * IRQL_DRIVER NO NO SORT OF YES (only higher priority drivers can be used)
 * IRQL_TIMER NO NO NO NO (but the timer handler jumps up to this level)
 * IRQL_HIGH NO NO NO NO
 *
 */

/*
 * Scheduler works. Page faults are allowed.
 */
#define IRQL_STANDARD 0
#define IRQL_STANDARD_HIGH_PRIORITY 1

/*
 * Scheduler still works at this point. Cannot page fault.
 */
#define IRQL_PAGE_FAULT 2

/*
 * This is the scheduler (and therefore things won't be scheduled out 'behind its back'). Cannot page fault.
 */
#define IRQL_SCHEDULER 3

/*
 * Scheduling will be postponed. Cannot page fault. Cannot use lower-priority devices.
 */
#define IRQL_DRIVER 4 // 3...39 is the driver range

/*
 * No scheduling, no page faulting, no using other hardware devices (no other irqs)
 */
#define IRQL_TIMER 40

/*
 * No interrupts from here.
 */
#define IRQL_HIGH 41

#include <common.h>
#include <assert.h>

void PostponeScheduleUntilStandardIrql(void);
void DeferUntilIrql(int irql, void(*handler)(void*), void* context);
int GetIrql(void);
int RaiseIrql(int level);
void LowerIrql(int level);
int GetNumberInDeferQueue(void);

void InitIrql(void);

#define MAX_IRQL(l) assert(GetIrql() <= l)
#define MIN_IRQL(l) assert(GetIrql() >= l)
#define EXACT_IRQL(l) assert(GetIrql() == l)
```

File: ./include/vnode.h

```
#pragma once

#include <common.h>
#include <sys/types.h>
#include <transfer.h>
#include <spinlock.h>
#include <sys/stat.h>

struct vnode;

/*
 * Operations which can be performed on an abstract file. They may be left NULL,
 * in this case, a default return value is supplied.
 *
 * check_open: default 0
 * Called just before a file is opened to ensure that the flags and the filename
 * are valid. Flags that can be passed in are O_RDONLY, O_WRONLY and O_RDWR, and
 * O_NONBLOCK. A filename may be invalid if the name is too long for the filesystem,
 * or if the filesystem contains other reserved characters.
 *
 * read: default EINVAL
 * Reads data from the file. If the file gives DT_DIR when asked for
 * dirent_type, then it should read in chunks of sizeof(struct dirent),
 * with the last being full of null bytes.
 *
 * write: default EINVAL
 * Writes data to the file. Fails on directories (EISDIR).
 *
 * ioctl: default EINVAL
```

```

* Performs a miscellaneous operation on a file.
*
* close: default 0
* Frees the vnode, as its reference count has hit zero.
*
* truncate: default EINVAL
* Truncates the file to the given size. Fails on directories (EISDIR).
*
* create: default EINVAL
* Creates a new file under a given parent, with a given name.
* The flags specifies O_RDWR, O_RDONLY, O_WRONLY, O_EXCL and O_APPEND.
*
* follow: default ENOTDIR
* Returns the vnode associated with a child of the current vnode.
* Fails on files (ENOTDIR).
*/

#define VNODE_WAIT_READ (1 << 0)
#define VNODE_WAIT_WRITE (1 << 1)
#define VNODE_WAIT_ERROR (1 << 2)
#define VNODE_WAIT_HAVE_TIMEOUT (1 << 3)
#define VNODE_WAIT_NON_BLOCK (1 << 4)

// TODO: the 'stat' data should live within the vnode... and then the stat
// call just returns it... file operations can then just adjust the stat struct
// as needed (e.g. when changing a file size in write).

struct vnode_operations {
int (*check_open)(struct vnode* node, const char* name, int flags);
int (*read)(struct vnode* node, struct transfer* io);
int (*write)(struct vnode* node, struct transfer* io);
int (*ioctl)(struct vnode* node, int command, void* buffer);
int (*close)(struct vnode* node); // release the filesystem specific data
int (*truncate)(struct vnode* node, off_t offset);
int (*create)(struct vnode* node, struct vnode** out, const char* name, int flags, mode_t mode);
int (*follow)(struct vnode* node, struct vnode** out, const char* name);
int (*wait)(struct vnode* node, int flags, uint64_t timeout_ms);

/*
* Must fail with EISDIR on directories. Should only decrement st.st_nlink,
* and remove the link from the filesystem. On things like FAT, where hard
* links are not supported, this can just decrement st.st_nlink, as we know
* that ops.delete is on its way, and that can properly delete it.
*/
int (*unlink)(struct vnode* node);

/*
* Deletes a file or directory from the filesystem completely. For files,
* the return value given will not propagate back to the VFS caller, as it
* gets called in DestroyVnode(). For files, st.st_nlink will be 0 on
* call.
*
* For directories, this function *must* check if the directory is non-empty
* and fail with ENOTEMPTY if so. st.st_nlink will be 1 on call - does not
* need to be modified.
*/
int (*delete)(struct vnode* node);
};

struct vnode {
struct vnode_operations ops;
void* data;
int reference_count;
struct spinlock reference_count_lock;
struct stat stat;
};

/*
* Allocates a new vnode for a given set of operations.
*/
struct vnode* CreateVnode(struct vnode_operations ops, struct stat st);
void ReferenceVnode(struct vnode* node);
void DereferenceVnode(struct vnode* node);

/*
* Wrapper functions to check the vnode is valid, and then call the driver.
*/
int VnodeOpCheckOpen(struct vnode* node, const char* name, int flags);
int VnodeOpRead(struct vnode* node, struct transfer* io);
int VnodeOpWrite(struct vnode* node, struct transfer* io);
int VnodeOpIoctl(struct vnode* node, int command, void* buffer);
int VnodeOpClose(struct vnode* node);
int VnodeOpTruncate(struct vnode* node, off_t offset);
uint8_t VnodeOpDirentType(struct vnode* node);
int VnodeOpCreate(struct vnode* node, struct vnode** out, const char* name, int flags, mode_t mode);
int VnodeOpFollow(struct vnode* node, struct vnode** out, const char* name);
int VnodeOpWait(struct vnode* node, int flags, uint64_t timeout_ms);
int VnodeOpUnlink(struct vnode* node);
int VnodeOpDelete(struct vnode* node);

```

File: `/include/blockingbuffer.h`

```

#pragma once

#include <common.h>

struct blocking_buffer;

struct blocking_buffer* BlockingBufferCreate(int size);
void BlockingBufferDestroy(struct blocking_buffer* buffer);
int BlockingBufferAdd(struct blocking_buffer* buffer, uint8_t c, bool block);
uint8_t BlockingBufferGet(struct blocking_buffer* buffer);
int BlockingBufferTryGet(struct blocking_buffer* buffer, uint8_t* c);

```

File: `/include/virtual.h`


```

#pragma once

#include <stddef.h>
#include <sys/types.h>
#include <stdbool.h>
#include <common.h>

#define VM_READ 1
#define VM_WRITE 2
#define VM_USER 4
#define VM_EXEC 8
#define VM_LOCK 16
#define VM_FILE 32
#define VM_FIXED_VIRT 64
#define VM_MAP_HARDWARE 128 /* map a physical page that doesn't live within the physical memory manager */
#define VM_MAP_LOCAL 256 /* indicates it's local to the VAS - i.e. not in kernel global memory */
#define VM_RECURSIVE 512 /* assumes the VAS is already locked, so won't lock or unlock it */
#define VM_RELOCATABLE 1024 /* needs driver fixups whenever swapped back in */
#define VM_EVICT_FIRST 2048

#define VAS_NO_ARCH_INIT 1

struct open_file;

struct vas_entry {
    size_t virtual;

    uint8_t in_ram : 1; /* Whether it is backed by a physical page or not. (i.e. does it have a real page table entry) */
    uint8_t allocated : 1; /* Whether or not to free a physical page on deallocation. Differs from in_ram when VM_MAP_HARDWARE is set. */
    uint8_t file : 1; /* Whether or not the page is file-mapped. */
    uint8_t cow : 1; /* */
    uint8_t swapfile : 1; /* Whether or not the page has been moved to a swapfile. Will not occur if 'file' is set (will back to that file instead) */
    uint8_t lock : 1;
    uint8_t read : 1;
    uint8_t write : 1;

    uint8_t exec : 1;
    uint8_t user : 1;
    uint8_t global : 1;
    uint8_t allow_temp_write : 1; /* used internally - allows the system to write to otherwise read-only pages to, e.g. reload from disk */
    uint8_t relocatable : 1; /* from a relocated driver file */
    uint8_t first_load : 1;
    uint8_t load_in_progress : 1; /* someone else is deferring a read into this page - keep trying the access until flag clears */

    uint8_t times_swapped : 4;
    uint8_t evict_first : 1;
    uint8_t : 3;

    int num_pages; /* only used for non-allocated or hardware mapped to reduce the number of AVL entries */

    off_t file_offset;
    struct open_file* file_node;
    size_t physical;
    union {
        size_t swapfile_offset;
        size_t relocation_base;
    };

    int ref_count;
};

struct vas;

size_t BytesToPages(size_t bytes);

bool LockVirt(size_t virtual);
void UnlockVirt(size_t virtual);
bool LockVirtEx(struct vas* vas, size_t virtual);
void UnlockVirtEx(struct vas* vas, size_t virtual);

void SetVirtPermissions(size_t virtual, int set, int clear);
int GetVirtPermissions(size_t virtual);
size_t MapVirt(size_t physical, size_t virtual, size_t bytes, int flags, struct open_file* file, off_t pos);
int UnmapVirt(size_t virtual, size_t bytes);
int UnmapVirtEx(struct vas* vas, size_t virtual, size_t pages);
size_t GetPhysFromVirt(size_t virtual);

struct vas* GetKernelVas(void); // a kernel vas
struct vas* GetVas(void); // current vas

struct vas* CreateVas(void);
void CreateVasEx(struct vas* vas, int flags);
void DestroyVas(struct vas* vas);

struct vas* CopyVas(void);
void SetVas(struct vas* vas);
void InitVirt(void);
bool IsVirtInitialised(void);
void EvictVirt(void);
void HandleVirtFault(size_t faulting_virt, int fault_type);

#include <arch.h>
#include <spinlock.h>

struct vas {
    struct avl_tree* mappings;
    platform_vas_data_t* arch_data;
    struct spinlock lock;
};

```

File: `/include/linkedlist.h`

```

#pragma once

#include <common.h>

```

```

struct linked_list;
struct linked_list_node;

struct linked_list* LinkedListCreate(void);
void LinkedListInsertStart(struct linked_list* list, void* data);
void LinkedListInsertEnd(struct linked_list* list, void* data);
bool LinkedListContains(struct linked_list* list, void* data);
int LinkedListGetIndex(struct linked_list* list, void* data);
void* LinkedListGetData(struct linked_list* list, int index);
bool LinkedListDeleteIndex(struct linked_list* list, int index);
bool LinkedListDeleteData(struct linked_list* list, void* data);
int LinkedListSize(struct linked_list* list);
void LinkedListDestroy(struct linked_list* list);

struct linked_list_node* LinkedListGetFirstNode(struct linked_list* list);
struct linked_list_node* LinkedListGetNextNode(struct linked_list_node* prev_node);
void* LinkedListGetDataFromNode(struct linked_list_node* node);

```

File: ./include/driver.h

```

#pragma once

#include <common.h>

struct quick_relocation {
    size_t address;
    size_t value;
};

struct quick_relocation_table {
    int total_entries;
    int used_entries;
    struct quick_relocation* entries;
};

struct vas;

void InitSymbolTable(void);
int RequireDriver(const char* name);
size_t GetDriverAddress(const char* name);
size_t GetSymbolAddress(const char* symbol);
void AddSymbol(const char* symbol, size_t address);

void SortQuickRelocationTable(struct quick_relocation_table* table);
void AddToQuickRelocationTable(struct quick_relocation_table* table, size_t addr, size_t val);
struct quick_relocation_table* CreateQuickRelocationTable(int count);
void PerformDriverRelocationOnPage(struct vas*, size_t relocation_base, size_t virt);

```

File: ./include/dev.h

```

#pragma once

void InitNullDevice(void);
void InitRandomDevice(void);

```

File: ./include/vfs.h

```

#pragma once

#include <common.h>
#include <sys/types.h>
#include <openfile.h>
#include <vnode.h>
#include <transfer.h>

void InitVfs(void);
int AddVfsMount(struct vnode* node, const char* name);
int RemoveVfsMount(const char* name);

int OpenFile(const char* path, int flags, mode_t mode, struct open_file** out);
int ReadFile(struct open_file* file, struct transfer* io);
int WriteFile(struct open_file* file, struct transfer* io);
int CloseFile(struct open_file* file);
int RemoveFileOrDirectory(const char* path, bool rmdir);

```

File: ./include/diskcache.h

```

#pragma once

#include <common.h>

#define DISKCACHE_NORMAL 0
#define DISKCACHE_REDUCE 1
#define DISKCACHE_TOSS 2

void InitDiskCaches(void);
void SetDiskCaches(int mode);

struct open_file* CreateDiskCache(struct open_file* underlying_disk);

```

File: ./include/timer.h

```

#pragma once

#include <common.h>

export uint64_t GetSystemTimer(void);

void ReceivedTimer(uint64_t nanos);
void InitTimer(void);

```

```

/*
 * Internal functions to do shenanigans
 */
struct thread;
void QueueForSleep(struct thread* thr);
bool TryDequeueForSleep(struct thread* thr);

```

File: ./include/unicode.h

```

#pragma once

#include <common.h>

int Utf16ToCodepoints(uint16_t* utf16, uint32_t* codepoints, int in_length, int* out_length);
int Utf8ToCodepoints(uint8_t* utf8, uint32_t* codepoints, int in_length, int* out_length);

int CodepointsToUtf16(uint32_t* codepoints, uint16_t* utf16, int in_length, int* out_length);
int CodepointsToUtf8(uint32_t* codepoints, uint8_t* utf8, int in_length, int* out_length);

```

File: ./include/swapfile.h

```

#pragma once

#include <common.h>

void InitSwapfile(void);

struct open_file* GetSwapfile(void);
uint64_t AllocateSwapfileIndex(void);
void DeallocateSwapfileIndex(uint64_t index);
int GetNumberOfPagesOnSwapfile(void);

```

File: ./include/voidptr.h

```

#pragma once
#include <stdint.h>

#define AddVoidPtr(ptr, offset) ((void*) (((uint8_t*) ptr) + offset))
#define SubVoidPtr(ptr, offset) ((void*) (((uint8_t*) ptr) - offset))

```

File: ./include/thread.h

```

#pragma once

#include <common.h>

struct semaphore;
struct process;

#define THREAD_STATE_RUNNING 0
#define THREAD_STATE_READY 1
#define THREAD_STATE_SLEEPING 2
#define THREAD_STATE_WAITING_FOR_SEMAPHORE 3
#define THREAD_STATE_WAITING_FOR_SEMAPHORE_WITH_TIMEOUT 4
#define THREAD_STATE_TERMINATED 5

#define SCHEDULE_POLICY_FIXED 0
#define SCHEDULE_POLICY_USER_HIGHER 1
#define SCHEDULE_POLICY_USER_NORMAL 2
#define SCHEDULE_POLICY_USER_LOWER 3

#define FIXED_PRIORITY_KERNEL_HIGH 0
#define FIXED_PRIORITY_KERNEL_NORMAL 30
#define FIXED_PRIORITY_IDLE 255

/*
 * Determines which of the 'next' pointers are used to manage the list.
 * A thread can be on multiple lists so long as they are different numbers.
 * Can increase the number of 'next' pointers in the thread struct to make them distinct if needed.
 */
#define NEXT_INDEX_READY 0
#define NEXT_INDEX_SLEEP 1
#define NEXT_INDEX_SEMAPHORE 2
#define NEXT_INDEX_TERMINATED 0 // terminated can share the ready list

struct thread {
/*
 * These first two values must be in this order.
 */
size_t kernel_stack_top;
size_t stack_pointer;

struct vas* vas;
size_t kernel_stack_size;
void (*initial_address)(void*);

/*
 * Allows a thread to be on a timer and a semaphore list at the same time.
 * Very sketchy stuff.
 */
struct thread* next[3];

int thread_id;
int state;
void* argument;
uint64_t time_used;
char* name;
int priority;
int schedule_policy;
size_t canary_position;
bool timed_out;
bool needs_termination;

```

```

struct semaphore* waiting_on_semaphore;

struct process* process;

/*
 * The system time at which this task's time has expired. If this is 0, then the task will not have a set time limit.
 * This value is set to GetSystemTimer() + TIMESLICE_LENGTH_MS when the task is scheduled in, and doesn't change until
 * the next time it is switched in.
 */
uint64_t timeslice_expiry;

uint64_t sleep_expiry;
};

void Schedule(void);
void LockSchedulerX(void);
void UnlockSchedulerX(void);
#define LockScheduler() /*LogWriteSerial("LOCKING SCHEDULER: %s, %s, %d\n", __FILE__, __func__, __LINE__);*/ LockSchedulerX()
#define UnlockScheduler() /*LogWriteSerial("UNLOCKING SCHEDULER: %s, %s, %d\n", __FILE__, __func__, __LINE__);*/ UnlockSchedulerX()

void InitScheduler(void);
void StartMultitasking(void);

void AssertSchedulerLockHeld(void);

struct thread* GetThread(void);
void TerminateThread(struct thread* thr);
void TerminateThreadLockHeld(struct thread* thr);

struct thread* CreateThreadEx(void(*entry_point)(void*), void* argument, struct vas* vas, const char* name, struct process* prcss, int policy, int priority, int kernel_stack_kb);
struct thread* CreateThread(void(*entry_point)(void*), void* argument, struct vas* vas, const char* name);

void BlockThread(int reason);
void UnblockThread(struct thread* thr);
int SetThreadPriority(struct thread* thread, int policy, int priority);

void SleepUntil(uint64_t system_time_ns);
void SleepNano(uint64_t delta_ns);
void SleepMilli(uint32_t delta_ms);

void HandleSleepWakeups(void* sys_time_ptr); // used internally between timer.c and thread.c
void InitIdle(void);
void InitCleaner(void);

struct process* CreateUserModeProcess(struct process* parent, const char* filename);

/*
 * A thread can lock itself onto the current cpu. Task switches *STILL OCCUR*, but we ensure that
 * next time this task runs, it will go back to this cpu.
 *
 * This is not a spinlock nor mutex, it's literally should just set a flag in the thread struct (sure, that
 * will spin while setting variable, but that's it). Between AssignThreadToCpu and UnassignThreadToCpu we remain
 * at IRQL_STANDARD.
 */
void AssignThreadToCpu(void);
void UnassignThreadToCpu(void);

```

File: `/include/common.h`

```

#pragma once

#include <stdint.h>
#include <stddef.h>
#include <stdbool.h>
#include <stdarg.h>

#define export __attribute__((used))

#ifndef NULL
#define NULL ((void*) 0)
#endif

#define warn_unused __attribute__((warn_unused_result))
#define always_inline __attribute__((always_inline)) inline

#define PAGEABLE_CODE_SECTION __attribute__((section("__pageablektext")))
#define PAGEABLE_DATA_SECTION __attribute__((section("__pageablekdata")))

#define NO_EXPORT __attribute__((visibility("hidden")))
#define EXPORT __attribute__((visibility("default")))

#define LOCKED_DRIVER_CODE __attribute__((section("__lockedtext")))
#define LOCKED_DRIVER_DATA __attribute__((section("__lockeddata")))
#define LOCKED_DRIVER_RODATA __attribute__((section("__lockedrodata")))

#define inline_memcpy(dst, src, n) __builtin_memcpy(dst, src, n)
#define inline_memset(dst, v, n) __builtin_memset(dst, v, n)

#define MAX(a, b) ((a) > (b) ? (a) : (b))
#define MIN(a, b) ((a) < (b) ? (a) : (b))
#define CLAMP(val, min, max) MAX(MIN(val, max), min)
#define COMPARE_SIGN(a, b) ((a) > (b) ? 1 : ((a) < (b) ? -1 : 0))

```

File: `/include/priorityqueue.h`

```

#pragma once

#include <stdbool.h>

struct priority_queue;

struct priority_queue_result {
    uint64_t priority;
    void* data;
};

```

```

};

struct priority_queue* PriorityQueueCreate(int capacity, bool max, int element_width);
void PriorityQueueInsert(struct priority_queue* queue, void* elem, uint64_t priority);
struct priority_queue_result PriorityQueuePeek(struct priority_queue* queue);
void PriorityQueuePop(struct priority_queue* queue);
int PriorityQueueGetCapacity(struct priority_queue* queue);
int PriorityQueueGetUsedSize(struct priority_queue* queue);
void PriorityQueueDestroy(struct priority_queue* queue);

```

File: ./include/avl.h

```

#pragma once

#include <common.h>

struct avl_tree;
struct avl_node;

typedef void (*avl_deletion_handler)(void*);
typedef int (*avl_comparator)(void*, void*);

void AvlTreePrint(struct avl_tree* tree, void(*printer)(void*));
struct avl_tree* AvlTreeCreate(void);
void AvlTreeInsert(struct avl_tree* tree, void* data);
void AvlTreeDelete(struct avl_tree* tree, void* data);
bool AvlTreeContains(struct avl_tree* tree, void* data);
void* AvlTreeGet(struct avl_tree* tree, void* data);
int AvlTreeSize(struct avl_tree* tree);
void AvlTreeDestroy(struct avl_tree* tree);
struct avl_node* AvlTreeGetRootNode(struct avl_tree* tree);
struct avl_node* AvlTreeGetLeft(struct avl_node* node);
struct avl_node* AvlTreeGetRight(struct avl_node* node);
void* AvlTreeGetData(struct avl_node* node);
avl_deletion_handler AvlTreeSetDeletionHandler(struct avl_tree* tree, avl_deletion_handler handler);
avl_comparator AvlTreeSetComparator(struct avl_tree* tree, avl_comparator comparator);

```

File: ./include/openfile.h

```

#pragma once

#include <common.h>
#include <sys/types.h>
#include <spinlock.h>

struct vnode;

struct open_file {
    bool can_read;
    bool can_write;
    mode_t initial_mode;
    size_t seek_position;
    int flags;
    int reference_count;
    struct spinlock reference_count_lock;
    struct vnode* node;
};

struct open_file* CreateOpenFile(struct vnode* node, int mode, int flags, bool can_read, bool can_write);
void ReferenceOpenFile(struct open_file* file);
void DereferenceOpenFile(struct open_file* file);

```

File: ./include/log.h

```

#pragma once

#include <common.h>

export void LogWriteSerial(const char* format, ...);
export void LogDeveloperWarning(const char* format, ...);

void DbgScreenPrintf(const char* format, ...);
void DbgScreenPuts(char* str);
void DbgScreenPutchar(char c);

```

File: ./include/threadlist.h

```

#pragma once

#include <common.h>

struct thread;

struct thread_list {
    struct thread* head;
    struct thread* tail;
    int index;
};

void ThreadListInit(struct thread_list* list, int index);
void ThreadListInsert(struct thread_list* list, struct thread* thread);
bool ThreadListContains(struct thread_list* list, struct thread* thread);
void ThreadListDelete(struct thread_list* list, struct thread* thread);
struct thread* ThreadListDeleteTop(struct thread_list* list);

```

File: ./include/partition.h

```

#pragma once

#include <stdint.h>
#include <stddef.h>

```

```

struct open_file;
struct vnode;

struct open_file* CreatePartition(struct open_file* disk, uint64_t start, uint64_t length, int id, int sector_size, int media_type, bool boot);
struct open_file** GetPartitionsForDisk(struct open_file* disk);

```

File: `/include/stdckdint.h`

```

#if __has_include(<stdckdint.h>)
# include <stdckdint.h>
#else
# ifdef __GNUC__
# define ckd_add(R, A, B) __builtin_add_overflow((A), (B), (R))
# define ckd_sub(R, A, B) __builtin_sub_overflow((A), (B), (R))
# define ckd_mul(R, A, B) __builtin_mul_overflow((A), (B), (R))
# else
# error "we need a compiler extension for this"
# endif
#endif

```

File: `/include/cpu.h`

```

#pragma once

#include <common.h>
#include <arch.h>
#include <spinlock.h>

struct vas;
struct thread;
struct avl_tree;

struct cpu {
    struct vas* current_vas;
    struct thread* current_thread;
    platform_cpu_data_t* platform_specific;
    size_t cpu_number;
    int irq;

    struct priority_queue* deferred_functions;
    struct priority_queue* irq_deferred_functions;
    bool init_irq_done;
    bool postponed_task_switch;

    struct avl_tree* global_vas_mappings;
    struct spinlock global_mappings_lock;
};

void InitCpuTable(void);
void InitBootstrapCpu(void);
void InitOtherCpu(void);

struct cpu* GetCpu(void);
int GetCpuCount(void);
struct cpu* GetCpuAtIndex(int index);

```

File: `/include/video.h`

```

#pragma once

struct video_driver {
    void (*putchar)(char);
    void (*puts)(char*);
};

void InitVideoConsole(struct video_driver driver);

```

File: `/include/heap.h`

```

#pragma once

#include <common.h>

/*
 * Allocation will not fault. If this is impossible to achieve, the system will panic.
 * Must not be set with HEAP_ALLOW_PAGING.
 */
#define HEAP_NO_FAULT 1

/*
 * Clears allocated memory to zero.
 */
#define HEAP_ZERO 2

/*
 * Indicates that the allocated region is allowed to be swapped onto disk.
 * Must not be set with HEAP_NO_FAULT. Data allocated with this flag set can only be
 * accessed when IRQL = IRQL_STANDARD.
 */
#define HEAP_ALLOW_PAGING 4

void* AllocHeap(size_t size);
void* AllocHeapEx(size_t size, int flags);
void* ReallocHeap(void* ptr, size_t size);
void* AllocHeapZero(size_t size);
void FreeHeap(void* ptr);
void ReinitHeap(void);
void InitHeap(void);

#ifdef NDEBUG
int DbgGetOutstandingHeapAllocations(void);

```

```
#endif
```

```
#define malloc(x) AllocHeap(x)
#define free(x) FreeHeap(x)
```

File: `/include/irq.h`

```
#pragma once
```

```
#include <arch.h>
#include <common.h>
```

```
typedef int(irq_handler_t)(platform_irq_context_t*);
```

```
int RegisterIrqHandler(int irq_num, irq_handler_t handler);
void RespondToIrq(int irq_num, int required_irq1, platform_irq_context_t* context);
void UnhandledFault(void);
```

File: `/include/process.h`

```
#pragma once
```

```
#include <common.h>
#include <sys/types.h>
```

```
struct filesdes_table;
struct process;
```

```
void InitProcess(void);
struct process* CreateProcess(pid_t parent_pid);
struct process* ForkProcess(void);
pid_t WaitProcess(pid_t pid, int* status, int flags);
void KillProcess(int retv);
```

```
struct process* GetProcessFromPid(pid_t pid);
struct process* GetProcess(void);
pid_t GetPid(struct process* p);
```

```
struct filesdes_table* GetFileDescriptorTable(struct process* p);
```

```
void AddThreadToProcess(struct process* p, struct thread* thr);
struct process* CreateProcessWithEntryPoint(pid_t parent, void(*entry_point)(void*), void* arg);
```

File: `/include/transfer.h`

```
#pragma once
```

```
#include <common.h>
```

```
enum transfer_type {
    TRANSFER_INTRA_KERNEL,
    TRANSFER_USERMODE,
};
```

```
enum transfer_direction {
    TRANSFER_READ,
    TRANSFER_WRITE,
};
```

```
/*
 * A data structure for performing file read and write operations, potentially
 * between the kernel and the user.
 *
 * TODO: userspace handling
 */
```

```
struct transfer {
    void* address;
    uint64_t length_remaining; /* In bytes. Will be modified on copying */
    uint64_t offset; /* In bytes. Will be modified on copying */
};
```

```
enum transfer_direction direction;
enum transfer_type type;
};
```

```
int PerformTransfer(void* trusted_buffer, struct transfer* untrusted_buffer, uint64_t len);
```

```
/*
 * max_length of 0 means unbounded
 */
```

```
int WriteStringToUsermode(const char* trusted_string, char* untrusted_buffer, uint64_t max_length);
int ReadStringFromUsermode(char* trusted_buffer, const char* untrusted_string, uint64_t max_length);
```

```
int WriteWordToUsermode(size_t* location, size_t value);
int ReadWordFromUsermode(size_t* location, size_t* output);
```

```
struct transfer CreateKernelTransfer(void* addr, uint64_t length, uint64_t offset, int direction);
struct transfer CreateTransferWritingToUser(void* untrusted_addr, uint64_t length, uint64_t offset);
struct transfer CreateTransferReadingFromUser(const void* untrusted_addr, uint64_t length, uint64_t offset);
```

File: `/include/panic.h`

```
#pragma once
```

```
#include <common.h>
```

```
enum {
    PANIC_UNKNOWN,
};
```

```
/*
 * A non-returnable function or infinite loop was exited out of.
 */
```

```

*/
PANIC_IMPOSSIBLE_RETURN,

/*
 * The panic was requested by the debugger.
 */
PANIC_MANUALLY_INITIATED,

/*
 * A unit test succeeded. Only to be used with the unit testing framework,
 * which panics to either succeed or fail (via assertion fails).
 */
PANIC_UNIT_TEST_OK,

/*
 * Used by drivers to report unrecoverable faults.
 */
PANIC_DRIVER_FAULT,

/*
 * The kernel heap is out of memory and cannot request any more.
 */
PANIC_OUT_OF_HEAP,

/*
 * Too much heap memory has been allocated before the virtual memory manager has
 * been initialised.
 */
PANIC_OUT_OF_BOOTSTRAP_HEAP,

/*
 * A request for a block on the heap was too large.
 */
PANIC_HEAP_REQUEST_TOO_LARGE,

/*
 * The kernel or a driver has caused an illegal page fault.
 */
PANIC_PAGE_FAULT_IN_NON_PAGED_AREA,

/*
 * An assertion failure within the kernel or driver.
 */
PANIC_ASSERTION_FAILURE,

/*
 * The bootloader failed to provide a usable memory map.
 */
PANIC_NO_MEMORY_MAP,

/*
 * The given section of kernel code is not implemented yet.
 */
PANIC_NOT_IMPLEMENTED,

/*
 * Wrong IRQL
 */
PANIC_INVALID_IRQL,

/*
 * Spinlock acquired from the wrong IRQL level.
 */
PANIC_SPINLOCK_WRONG_IRQL,

/*
 * No more physical memory, even after evicting old pages.
 */
PANIC_OUT_OF_PHYS,

PANIC_PRIORITY_QUEUE,
PANIC_LINKED_LIST,

/*
 * Kernel stack overflow
 */
PANIC_CANARY_DIED,

PANIC_SEMAPHORE_DESTROY_WHILE_HELD,
PANIC_SEM_BLOCK_WITHOUT_THREAD,
PANIC_CANNOT_LOCK_MEMORY,
PANIC_THREAD_LIST,
PANIC_CANNOT_MALLOC_WITHOUT_FAULTING,
PANIC_NO_FILESYSTEM,
PANIC_BAD_KERNEL,
PANIC_DISK_FAILURE_ON_SWAPFILE,
PANIC_NEGATIVE_SEMAPHORE,
PANIC_NON_MASKABLE_INTERRUPT,
PANIC_UNHANDLED_KERNEL_EXCEPTION,
PANIC_REQUIRED_DRIVER_MISSING_SYMBOL,
PANIC_REQUIRED_DRIVER_NOT_FOUND,
PANIC_NO_LOW_MEMORY,
PANIC_OUT_OF_SWAPFILE,
PANIC_PROGRAM_LOADER,
PANIC_VAS_TRIED_TO_SELF_DESTRUCT,
PANIC_ACPI_AML,
PANIC_SPINLOCK_DOUBLE_ACQUISITION,
PANIC_SPINLOCK_RELEASED_BEFORE_ACQUIRED,
PANIC_DOUBLE_FREE_DETECTED,
PANIC_CONFLICTING_ALLOCATION_REQUIREMENTS,

_PANIC_HIGHEST_VALUE
};

_Noreturn void PanicEx(int code, const char* message);
_Noreturn void Panic(int code);

```



```
const char* GetPanicMessageFromCode(int code);
int SetGraphicalPanicHandler(void (*handler)(int, const char*));
```

File: `/include/diskutil.h`

```
#pragma once

#include <common.h>

struct vnode;
struct open_file;

#define MAX_PARTITIONS_PER_DISK 8

#define DISKUTIL_TYPE_FIXED 0
#define DISKUTIL_TYPE_FLOPPY 1
#define DISKUTIL_TYPE_OPTICAL 2
#define DISKUTIL_TYPE_REMOVABLE 3
#define DISKUTIL_TYPE_NETWORK 4
#define DISKUTIL_TYPE_VIRTUAL 5
#define DISKUTIL_TYPE_RAM 6
#define DISKUTIL_TYPE_OTHER 7

#define __DISKUTIL_NUM_TYPES 8

struct disk_partition_helper {
    struct vnode* partitions[MAX_PARTITIONS_PER_DISK];
    char* partition_names[MAX_PARTITIONS_PER_DISK];
    int num_partitions;
};

void InitDiskUtil(void);
char* GenerateNewRawDiskName(int type);
char* GenerateNewMountedDiskName();
void CreateDiskPartitions(struct open_file* disk);
void InitDiskPartitionHelper(struct disk_partition_helper* helper);

int DiskFollowHelper(struct disk_partition_helper* helper, struct vnode** out, const char* name);
int DiskCreateHelper(struct disk_partition_helper* helper, struct vnode** in, const char* name);
```

File: `/include/fs/internal/fat.h`

```
#pragma once

#include <common.h>

struct open_file;

#define LFN_SHORT_ONLY 0
#define LFN_BOTH 1
#define LFN_ERROR 2

#define FAT12 0
#define FAT16 2 // the value of 2 is relied on, as it means 2 bytes per FAT (is used for calcs)
#define FAT32 4 // as above, we use fat_type to do calcs, so required that FAT32 == 4

struct fat_data {
    int num_fats;
    int fat_sectors[4];
    int sectors_per_fat;
    union {
        uint64_t first_root_dir_sector_12_16;
        uint64_t root_dir_cluster_32;
    };
    uint64_t root_dir_num_sectors_12_16;
    int total_clusters;
    uint64_t first_data_sector;
    uint64_t first_fat_sector;
    int fat_type; // FAT12 or FAT16 or FAT32
    int sectors_per_cluster;
    int bytes_per_sector;
};

struct open_file* disk; // TODO! points to a vnode for the partition

uint8_t* cluster_buffer_a;
uint8_t* cluster_buffer_b;
};

int GetFatShortFilename(char* lfn, char* output, char* directory);
void FormatFatShortName(char* with_dot, char* without_dot);
void UnformatFatShortName(char* without_dot, char* with_dot);

int ReadFatCluster(struct fat_data* fat, int cluster, bool buffer);
int WriteFatCluster(struct fat_data* fat, int cluster, bool buffer);
int ReadFatEntry(struct fat_data* fat, int entry, uint32_t* output);
int WriteFatEntry(struct fat_data* fat, int entry, uint32_t value);

struct fat_data LoadFatData(uint8_t* boot_sector, struct open_file* disk);
```

File: `/include/fs/fat.h`

```
#pragma once

#include <common.h>

int DetectFatPartition(void* partition);
```

File: `/include/syscall.h`

```
#pragma once

#include <common.h>
```

```
int HandleSystemCall(int call, size_t a, size_t b, size_t c, size_t d, size_t e);
```

```
int SysYield(size_t, size_t, size_t, size_t, size_t);
int SysTerminate(size_t, size_t, size_t, size_t, size_t);
int SysMapVirt(size_t, size_t, size_t, size_t, size_t);
int SysUnmapVirt(size_t, size_t, size_t, size_t, size_t);
int SysOpen(size_t, size_t, size_t, size_t, size_t);
int SysReadWrite(size_t, size_t, size_t, size_t, size_t);
int SysClose(size_t, size_t, size_t, size_t, size_t);
int SysSeek(size_t, size_t, size_t, size_t, size_t);
int SysDup(size_t, size_t, size_t, size_t, size_t);
int SysExit(size_t, size_t, size_t, size_t, size_t);
```

File: `/include/physical.h`

```
#pragma once
```

```
#include <common.h>
```

```
void DeallocPhys(size_t addr);
void DeallocPhysContiguous(size_t addr, size_t bytes);
size_t AllocPhys(void);
size_t AllocPhysContiguous(size_t bytes, size_t min_addr, size_t max_addr, size_t boundary);
size_t GetTotalPhysKilobytes(void);
size_t GetFreePhysKilobytes(void);
```

```
void InitPhys(void);
void ReinitPhys(void);
```

File: `/include/debug/tfw_tests.h`

```
#pragma once
```

```
#ifndef NDEBUG
```

```
void RegisterTfwPhysTests(void);
void RegisterTfwInitTests(void);
void RegisterTfwIrqTests(void);
void RegisterTfwAVLTreeTests(void);
void RegisterTfwPriorityQueueTests(void);
void RegisterTfwSemaphoreTests(void);
void RegisterTfwWaitTests(void);
```

```
#endif
```

File: `/include/debug/tfw.h`

```
#pragma once
```

```
#include <common.h>
```

```
enum {
    TFW_SP_INITIAL,
    TFW_SP_AFTER_PHYS,
    TFW_SP_AFTER_HEAP,
    TFW_SP_AFTER_BOOTSTRAP_CPU,
    TFW_SP_AFTER_VIRT,
    TFW_SP_AFTER_PHYS_REINIT,
    TFW_SP_AFTER_ALL_CPU,

```

```
    TFW_SP_ALL_CLEAR,
};
```

```
#ifndef NDEBUG
```

```
#define IsInTfwTest() false
#define FinishedTfwTest(x)
#define MarkTfwStartPoint(x)
#define InitTfw()
```

```
#else
```

```
bool IsInTfwTest(void);
```

```
// this can probably go up to around 150,000 or so in theory (in what the transfer format supports), or about 20,000 on a 4MB RAM system.
```

```
// but bigger is slower, so only increase as we need to
```

```
#define MAX_TFW_TESTS 100
```

```
#define MAX_NAME_LENGTH 96 // If this changes the python must do too
```

```
struct tfw_test {
    char name[MAX_NAME_LENGTH];
    void (*code)(struct tfw_test*, size_t context);
    int start_point;
    int expected_panic_code;
    bool nightly_only;
    size_t context;
};
```

```
#define TFW_IGNORE_UNUSED (void) test; (void) context;
```

```
#define TFW_CREATE_TEST(name) static void name (struct tfw_test* test, size_t context)
```

```
void RegisterTfwTest(const char* name, int start_point, void (*code)(struct tfw_test*, size_t), int expected_panic, size_t context);
void RegisterNightlyTfwTest(const char* name, int start_point, void (*code)(struct tfw_test*, size_t), int expected_panic, size_t context);
```

```
void FinishedTfwTest(int panic_code);
void MarkTfwStartPoint(int id);
void InitTfw(void);
```

```
#endif
```

File: ./include/debug/hostio.h

```
#pragma once

#ifdef NDEBUG

#else

#include <common.h>

#define DBGPKT_TFW 0

void DbgWritePacket(int type, uint8_t* data, int size);
void DbgReadPacket(int* type, uint8_t* data, int* size);

#endif
```

File: ./include/arch.h

```
#pragma once

/*
 * arch.h - Architecture-specific wrappers
 *
 *
 * Functions relating to hardware devices that must be implemented by
 * any platform supporting the operating system.
 */

/*
 * config.h needs to define the following:
 * - ARCH_PAGE_SIZE
 * - ARCH_MAX_CPU_ALLOWED
 * - ARCH_MAX_RAM_KBS
 * - ARCH_BIG_ENDIAN or ARCH_LITTLE_ENDIAN
 * - the address in the kernel area, ARCH_PROG_LOADER_BASE, where the program loader lives, and
 * - ARCH_PROG_LOADER_ENTRY, the entry point of the prog loader
 * - the valid user area, via ARCH_USER_AREA_BASE and ARCH_USER_AREA_LIMIT
 * - the valid kernel area, via ARCH_KRNL_SBRK_BASE and ARCH_KRNL_SBRK_LIMIT
 * (the kernel and user areas must not overlap, but ARCH_USER_AREA_LIMIT may equal ARCH_KRNL_SBRK_BASE
 or ARCH_KRNL_SBRK_LIMIT may equal ARCH_USER_AREA_BASE)
 * - the user stack area, via ARCH_USER_STACK_BASE and ARCH_USER_STACK_LIMIT
 * (may overlap with ARCH_USER_AREA_BASE and ARCH_USER_AREA_LIMIT)
 * - a typedef for platform_cpu_data_t
 * - a typedef for platform_irq_context_t
 * - a typedef for platform_vas_data_t
 */

#include <machine/config.h>

#if ARCH_USER_STACK_BASE < ARCH_USER_AREA_BASE
#error "ARCH_USER_STACK_BASE must be greater than or equal to ARCH_USER_AREA_BASE"
#elif ARCH_USER_STACK_LIMIT > ARCH_USER_AREA_LIMIT
#error "ARCH_USER_STACK_LIMIT must be less than or equal to ARCH_USER_AREA_LIMIT"
#endif

#include <common.h>

struct arch_memory_range
{
    size_t start;
    size_t length;
};

struct vas;
struct vas_entry;
struct thread;
struct open_file;
struct cpu;
struct quick_relocation_table;

struct arch_driver_t;

/*
 * Only to be called in very specific places, e.g. turning interrupts
 * on for the first time, the panic handler.
 */
void ArchEnableInterrupts(void);
void ArchDisableInterrupts(void);

/*
 * Do nothing until (maybe) the next interrupt. If this is not supported by the
 * system it may just return without doing anything.
 */
void ArchStallProcessor(void);

#define ARCH_POWER_STATE_REBOOT 1
#define ARCH_POWER_STATE_SHUTDOWN 2
#define ARCH_POWER_STATE_SLEEP 3
int ArchSetPowerState(int power_state);

void ArchSpinlockAcquire(volatile size_t* lock);
void ArchSpinlockRelease(volatile size_t* lock);

/*
 * To be called repeatedly until it returns NULL. Each time will return a new memory
 * range. An address of a static local object is permitted to be returned.
 *
 * NULL is returned if there is no more memory. No more calls to this function
 * will be made after a NULL is returned.
 */
struct arch_memory_range* ArchGetMemory() warn_unused;

uint64_t ArchReadTimestamp(void);
```

```

void ArchFlushTlb(struct vas* vas);
void ArchAddMapping(struct vas* vas, struct vas_entry* entry);
void ArchUpdateMapping(struct vas* vas, struct vas_entry* entry);
void ArchUnmap(struct vas* vas, struct vas_entry* entry);
void ArchSetVas(struct vas* vas);

void ArchGetPageUsageBits(struct vas* vas, struct vas_entry* entry, bool* accessed, bool* dirty);
void ArchSetPageUsageBits(struct vas* vas, struct vas_entry* entry, bool accessed, bool dirty);

// responsible for loading all symbols. should not close the file!
int ArchLoadDriver(size_t* relocation_point, struct open_file* file, struct quick_relocation_table** table);
void ArchLoadSymbols(struct open_file* file, size_t adjust);
void ArchSwitchThread(struct thread* old, struct thread* new);
size_t ArchPrepareStack(size_t addr);

void ArchSwitchToUsermode(size_t entry_point, size_t user_stack, void* arg);

void ArchInitDev(bool fs);

/*
 * Used only if the AVL tree is insufficient, e.g. for deallocating part of the kernel region to, e.g.
 * reclaim the physical memory bitmap. Works only for the current VAS. Returns 0 on no mapping.
 */
size_t ArchVirtualToPhysical(size_t virtual);

/*
 * Initialises a given VAS with platform specific data (e.g. mapping the kernel in).
 */
void ArchInitVas(struct vas* vas);

/*
 * Initialises virtual memory in general, i.e. creates the first VAS.
 */
void ArchInitVirt(void);

int ArchGetCurrentCpuIndex(void);
void ArchSendEoi(int irq_num);
/*
 * Sets the CPUs interrupt state (and mask devices) based on an IRQL. This function
 * will always be called with interrupts completely disabled.
 */
void ArchSetIrql(int irql);

void ArchInitBootstrapCpu(struct cpu* cpu);

/*
 * If possible, initialises the next CPU, and returns true. If there are no more CPUs
 * to initialise, returns false.
 */
bool ArchInitNextCpu(struct cpu* cpu);

```

File: */irq/irq.c*

```

1
2 #include <arch.h>
3 #include <irq.h>
4 #include <log.h>
5 #include <irq.h>
6 #include <thread.h>
7 #include <linkedlist.h>
8 #include <errno.h>
9 #include <process.h>
10 #include <panic.h>
11
12 #define HIGHEST_IRQ_NUM 256
13
14 static struct linked_list* irq_table[HIGHEST_IRQ_NUM] = {0};
15
16 int RegisterIrqHandler(int irq_num, irq_handler_t handler) {
17     if (irq_num < 0 || irq_num >= HIGHEST_IRQ_NUM || handler == NULL) {
18         return EINVAL;
19     }
20
21     if (irq_table[irq_num] == NULL) {
22         irq_table[irq_num] = LinkedListCreate();
23     }
24
25     LinkedListInsertEnd(irq_table[irq_num], (void*)(size_t) handler);
26     return 0;
27 }
28
29 void RespondToIrq(int irq_num, int required_irq, platform_irq_context_t* context) {
30     int irq = RaiseIrq(required_irq);
31     ArchSendEoi(irq_num); /* this must be done after raising the IRQ */
32
33     if (irq_table[irq_num] != NULL) {
34         struct linked_list_node* iter = LinkedListGetFirstNode(irq_table[irq_num]);
35         while (iter != NULL) {
36             irq_handler_t handler = (irq_handler_t)(size_t) LinkedListGetDataFromNode(iter);
37             assert(handler != NULL);
38
39             /*
40              * Interrupt handlers return 0 if they could handle the IRQ (i.e. stop trying to handle it).
41              * Non-zero means 'leave this one for someone else'.
42              */
43             if (handler(context) == 0) {
44                 break;
45             }
46
47             iter = LinkedListGetNextNode(iter);
48         }
49     }
50
51     LowerIrq(irq);
52 }
53
54 void UnhandledFault(void) {
55     if (GetProcess() != NULL) {
56         LogWriteSerial("unhandled fault...\n");
57         TerminateThread(GetThread());
58     } else {
59         Panic(PANIC_UNHANDLED_KERNEL_EXCEPTION);
60     }
61 }
62

```

File: ./irq/cpu.c

```

1 #include <cpu.h>
2 #include <arch.h>
3 #include <assert.h>
4 #include <heap.h>
5 #include <string.h>
6 #include <spinlock.h>
7 #include <irq.h>
8
9 static struct cpu cpu_table[ARCH_MAX_CPU_ALLOWED];
10 static int num_cpus_running = 1;
11
12 static void InitCpuTableEntry(int index) {
13     /*
14      * The boot CPU can't use dynamic memory, as this happens before we have a heap.
15      */
16     static platform_cpu_data_t boot_cpu_data;
17
18     cpu_table[index].cpu_number = index;
19     cpu_table[index].platform_specific = index != 0 ? AllocHeapZero(sizeof(platform_cpu_data_t)) : &boot_cpu_data;
20     cpu_table[index].irq = IRQ_STANDARD;
21     cpu_table[index].global_vas_mappings = NULL;
22     cpu_table[index].current_vas = NULL;
23     cpu_table[index].current_thread = NULL;
24     cpu_table[index].init_irq_done = false;
25     cpu_table[index].postponed_task_switch = false;
26     InitSpinlock(&cpu_table[index].global_mappings_lock, "gbl vas map", IRQ_SCHEDULER);
27 }
28
29 /*
30  * Initialises the CPU table ('cpu_table') for the bootstrap processor. This does *not* do any
31  * platform-specific initialisation, as that requires other features to be set up first.
32  */
33 void InitCpuTable(void) {
34     assert(num_cpus_running == 1);
35     InitCpuTableEntry(0);
36 }
37
38 /*
39  * Performs platform-specific initialisation of the bootstrap CPU (e.g. initialising interrupts,
40  * system timers, segments, etc.).
41  */
42 void InitBootstrapCpu(void) {
43     EXACT_IRQ(IRQ_STANDARD);
44     ArchInitBootstrapCpu(cpu_table);
45 }
46
47 void InitOtherCpu(void) {
48     /*
49      * We initialise the next entry, even before we know there's a CPU there. If there isn't a
50      * CPU there, we don't increment 'num_cpus_running' so the entry won't get used.
51      */
52     InitCpuTableEntry(1);
53
54     while (ArchInitNextCpu(cpu_table + num_cpus_running)) {
55         ++num_cpus_running;
56         InitCpuTableEntry(num_cpus_running);
57     }
58 }
59
60 int GetCpuCount(void) {
61     return num_cpus_running;
62 }
63
64 struct cpu* GetCpuAtIndex(int index) {
65     assert(index >= 0 && index < GetCpuCount());
66     return cpu_table + index;
67 }
68
69 struct cpu* GetCpu(void) {
70     return cpu_table + ArchGetCurrentCpuIndex();
71 }

```

File: `/irq/irq.c`

```

1
2 #include <panic.h>
3 #include <cpu.h>
4 #include <irq.h>
5 #include <log.h>
6 #include <priorityqueue.h>
7 #include <thread.h>
8 #include <assert.h>
9
10 /**
11  * Runs a function at an IRQ level lower than or equal to the current IRQ level. If the IRQ levels match,
12  * the function will be run immediately. If the target IRQ level is lower than the current IRQ level,
13  * it will be deferred until the IRQ level drops below IRQ_LEVEL_SCHEDULER. Deferred function
14  * calls will be run in order from highest IRQ level to lowest. If InitIrql() has not been called
15  * prior to calling this function, any requests meant to be deferred will instead be silently
16  * ignored (this is needed to bootstrap the physical memory manager, et al.).
17  *
18  * @param irql The IRQ level to run the function at. This IRQ level must be lower than the
19  *             current IRQ level, or a panic will occur.
20  * @param handler The function to be run.
21  * @param context An argument given to the handler function.
22  */
23 void DeferUntilIrql(int irql, void(*handler)(void*), void* context) {
24     if (irql < GetIrql() || (irql == IRQ_LEVEL_STANDARD_HIGH_PRIORITY && GetIrql() == IRQ_LEVEL_STANDARD)) {
25         handler(context);
26     }
27     else if (irql > GetIrql()) {
28         PanicEx(PANIC_INVALID_IRQ, "invalid irql on DeferUntilIrql");
29     }
30     else if (GetCpu()->init_irql_done) {
31         struct irq_deferment = {.context = context, .handler = handler};
32         PriorityQueueInsert(GetCpu()->deferred_functions, (void*) &deferment, irql);
33     }
34 }
35
36 int GetIrql(void) {
37     return GetCpu()->irql;
38 }
39
40 int RaiseIrql(int level) {
41     ArchDisableInterrupts();
42
43     struct cpu* cpu = GetCpu();
44     int existing_level = cpu->irql;
45
46     if (level < existing_level) {
47         PanicEx(PANIC_INVALID_IRQ, "invalid irql on RaiseIrql");
48     }
49
50     cpu->irql = level;
51     ArchSetIrql(level);
52
53     return existing_level;
54 }
55
56 void LowerIrql(int target_level) {
57     struct cpu* cpu = GetCpu();
58     struct priority_queue* deferred_functions = cpu->deferred_functions;
59
60     int current_level = cpu->irql;
61
62     if (target_level > current_level) {
63         PanicEx(PANIC_INVALID_IRQ, "invalid irql on LowerIrql");
64     }
65
66     while (cpu->init_irql_done && PriorityQueueGetUsedSize(deferred_functions) > 0) {
67         struct priority_queue_result next = PriorityQueuePeek(deferred_functions);
68         assert(((int) next.priority <= current_level || (next.priority == IRQ_LEVEL_STANDARD_HIGH_PRIORITY && current_level == IRQ_LEVEL_STANDARD)));
69
70         if ((int) next.priority >= target_level) {
71             current_level = next.priority;
72             if (current_level == IRQ_LEVEL_STANDARD_HIGH_PRIORITY) {
73                 current_level = IRQ_LEVEL_STANDARD;
74             }
75
76             /*
77              * Must Pop() before we call the handler (otherwise if the handler does a raise/lower, it will
78              * retrigger itself and cause a recursion loop), and must also get data off the queue before we Pop().
79              * Also must only actually lower the IRQ level after doing this, so we don't get interrupted in between
80              * (as someone else could then Raise/Lower, and mess us up.)
81              */
82             struct irq_deferment* deferred_call = (struct irq_deferment*) next.data;
83             void* context = deferred_call->context;
84             void (*handler)(void*) = deferred_call->handler;
85             if (handler == NULL) {
86                 ArchSetIrql(current_level);
87                 continue;
88             }
89             PriorityQueuePop(deferred_functions);
90             cpu->irql = current_level;
91             ArchSetIrql(current_level);
92             handler(context);
93
94         } else {
95             break;
96         }
97     }
98
99     current_level = target_level;
100     if (current_level == IRQ_LEVEL_STANDARD_HIGH_PRIORITY) {
101         current_level = IRQ_LEVEL_STANDARD;
102     }
103     cpu->irql = current_level;
104     ArchSetIrql(current_level);
105
106     if (current_level == IRQ_LEVEL_STANDARD && cpu->postponed_task_switch) {
107         cpu->postponed_task_switch = false;
108         Schedule();
109     }
110 }
111
112 void PostponeScheduleUntilStandardIrql(void) {
113     // TODO: does this function need its own lock? (e.g. just for setting postponed_task_switch)
114     GetCpu()->postponed_task_switch = true;
115 }
116
117 /**
118  * Requires TFW_SP_AFTER_HEAP or later.
119  */
120 void InitIrql(void) {
121     GetCpu()->deferred_functions = PriorityQueueCreate(32, true, sizeof(struct irq_deferment));
122     GetCpu()->init_irql_done = true;
123 }
124
125 int GetNumberInDeferQueue(void) {
126     return PriorityQueueGetUsedSize(GetCpu()->deferred_functions);
127 }

```

File: ./sys/calls/exit.c

```
1 #include <syscall.h>
2 #include <errno.h>
3 #include <_syscallnum.h>
4 #include <thread.h>
5 #include <process.h>
6
7 int SysExit(size_t status, size_t, size_t, size_t, size_t) {
8     KillProcess(status);
9     return EFAULT;
10 }
```

File: ./sys/calls/mapvirt.c

```
1 #include <syscall.h>
2 #include <errno.h>
3 #include <_syscallnum.h>
4 #include <thread.h>
5 #include <transfer.h>
6 #include <process.h>
7 #include <filedes.h>
8 #include <virtual.h>
9 #include <stdckdint.h>
10 #include <sys/mman.h>
11
12 int SysMapVirt(size_t flags, size_t bytes, size_t fd, size_t offset, size_t virtual) {
13     size_t* userptr_virt = (size_t*)virtual;
14
15     if (flags & ~(VM_READ | VM_WRITE | VM_EXEC | VM_FILE | VM_FIXED_VIRT)) {
16         return EINVAL;
17     }
18
19     size_t target_virtual;
20     int res = ReadWordFromUsermode(userptr_virt, &target_virtual);
21     if (res != 0) {
22         return res;
23     }
24
25     if (target_virtual < ARCH_USER_AREA_BASE) {
26         return EINVAL;
27     }
28
29     size_t end_of_virtual;
30     bool overflow = ckd_add(&end_of_virtual, target_virtual, bytes);
31
32     if (overflow || (end_of_virtual >= ARCH_USER_AREA_LIMIT)) {
33         return EINVAL;
34     }
35
36     struct open_file* file = NULL;
37     if (flags & VM_FILE) {
38         res = GetFileFromDescriptor(GetFileDescriptorTable(GetProcess()), fd, &file);
39         if (file == NULL || res != 0) {
40             return res;
41         }
42     }
43
44     size_t output_virtual = MapVirt(0, target_virtual, bytes, flags | VM_USER | VM_LOCAL, file, offset);
45     if (output_virtual == 0) {
46         return EINVAL;
47     }
48
49     return WriteWordToUsermode(userptr_virt, output_virtual);
50 }
```

File: ./sys/calls/seek.c

```
1 #include <syscall.h>
2 #include <errno.h>
3 #include <_syscallnum.h>
4 #include <thread.h>
5 #include <log.h>
6 #include <vfs.h>
7 #include <process.h>
8 #include <dirent.h>
9 #include <filedes.h>
10 #include <unistd.h>
11 #include <transfer.h>
12
13 int SysSeek(size_t fd, size_t pos_ptr, size_t whence, size_t, size_t) {
14     struct open_file* file;
15     int res = GetFileFromDescriptor(GetFileDescriptorTable(GetProcess()), fd, &file);
16
17     if (file == NULL || res != 0) {
18         return res;
19     }
20
21     struct transfer_io = CreateTransferReadingFromUser((void*) pos_ptr, sizeof(off_t), 0);
22     off_t offset;
23     if ((res = PerformTransfer(&offset, &io, sizeof(off_t))) {
24         return res;
25     }
26
27     int type = VnodeOpDirentType(file->node);
28     if (type == DT_FIFO || type == DT_SOCKET) {
29         return ESPIPE;
30     }
31
32     if (whence == SEEK_CUR) {
33         offset += file->seek_position;
34     }
35     else if (whence == SEEK_END) {
36         offset += file->node->stat.st_size;
37     }
38     else if (whence != SEEK_SET) {
39         return EINVAL;
40     }
41
42     file->seek_position = offset;
43
44     io = CreateTransferWritingToUser((void*) pos_ptr, sizeof(off_t), 0);
45     return PerformTransfer(&offset, &io, sizeof(off_t));
46 }
```

File: ./sys/calls/unmapvirt.c


```

1 #include <syscall.h>
2 #include <errno.h>
3 #include <_syscallnum.h>
4 #include <thread.h>
5 #include <transfer.h>
6 #include <virtual.h>
7 #include <stdckdint.h>
8 #include <sys/mman.h>
9
10 int SysUnmapVirt(size_t virtual, size_t bytes, size_t, size_t, size_t) {
11     if (virtual < ARCH_USER_AREA_BASE) {
12         return EINVAL;
13     }
14
15     size_t end_of_virtual;
16     bool overflow = ckd_add(&end_of_virtual, virtual, bytes);
17
18     if (overflow || (end_of_virtual >= ARCH_USER_AREA_LIMIT)) {
19         return EINVAL;
20     }
21
22     return UnmapVirt(virtual, bytes);
23 }

```

File: ./sys/calls/readwrite.c

```

1 #include <syscall.h>
2 #include <errno.h>
3 #include <_syscallnum.h>
4 #include <thread.h>
5 #include <log.h>
6 #include <vfs.h>
7 #include <transfer.h>
8 #include <fontl.h>
9 #include <process.h>
10 #include <filedes.h>
11
12 int SysReadWrite(size_t fd, size_t size, size_t buffer, size_t br_out, size_t write) {
13     struct filedes_table* table = GetFileDescriptorTable(GetProcess());
14     struct open_file* file;
15     int res = GetFileFromDescriptor(table, fd, &file);
16
17     if (file == NULL || res != 0) {
18         return res;
19     }
20
21     if (write && (file->flags & O_APPEND)) {
22         file->seek_position = file->node->stat.st_size;
23     }
24
25     struct transfer io = CreateTransferWritingToUser((uint8_t*) buffer, size, file->seek_position);
26     io.direction = write ? TRANSFER_WRITE : TRANSFER_READ;
27
28     if ((res = (write ? WriteFile : ReadFile)(file, &io))) {
29         return res;
30     }
31
32     size_t br = size - io.length_remaining;
33     file->seek_position += br;
34
35     return WriteWordToUsermode((size_t*) br_out, br);
36 }

```

File: ./sys/calls/terminate.c

```

1 #include <syscall.h>
2 #include <errno.h>
3 #include <_syscallnum.h>
4 #include <thread.h>
5 #include <process.h>
6 #include <assert.h>
7
8 int SysTerminate(size_t status, size_t, size_t, size_t, size_t) {
9     assert(GetProcess() != NULL);
10
11     // TODO: check if last thread in process
12     // if (GetNumberOfRemainingThreads(GetProcess()) == 1) {
13     //     KillProcess(status);
14     // } else {
15     //     TerminateThread(GetThread());
16     // }
17
18     (void) status;
19
20     TerminateThread(GetThread());
21     return EFAULT;
22 }

```

File: ./sys/calls/dup.c

```

1 #include <syscall.h>
2 #include <errno.h>
3 #include <_syscallnum.h>
4 #include <thread.h>
5 #include <log.h>
6 #include <vfs.h>
7 #include <process.h>
8 #include <filedes.h>
9 #include <fcntl.h>
10
11 int SysDup(size_t dup_num, size_t old_fd, size_t new_fd, size_t flags, size_t) {
12     struct filedes_table* table = GetFileDescriptorTable(GetProcess());
13
14     if ((flags & ~O_CLOEXEC) != 0) {
15         return EINVAL;
16     }
17
18     if (dup_num == 1) {
19         int result_fd;
20         int res = DuplicateFileDescriptor(table, old_fd, &result_fd);
21         if (res != 0) {
22             return res;
23         }
24
25         return WriteWordToUsermode((size_t*) new_fd, result_fd);
26     } else if (dup_num == 2) {
27         return DuplicateFileDescriptor2(table, old_fd, new_fd, flags);
28     } else {
29         return EINVAL;
30     }
31 }
32
33

```

File: ./sys/calls/yield.c

```

1 #include <syscall.h>
2 #include <errno.h>
3 #include <_syscallnum.h>
4 #include <thread.h>
5 #include <log.h>
6
7 int SysYield(size_t, size_t, size_t, size_t, size_t) {
8     Schedule();
9     return 0;
10 }

```

File: ./sys/calls/open.c

```

1 #include <syscall.h>
2 #include <errno.h>
3 #include <_syscallnum.h>
4 #include <thread.h>
5 #include <transfer.h>
6 #include <log.h>
7 #include <fcntl.h>
8 #include <thread.h>
9 #include <process.h>
10 #include <vfs.h>
11 #include <filedes.h>
12
13 #define ALLOWABLE_FLAGS (O_CREAT | O_EXCL | O_NOCTTY | O_TRUNC | O_APPEND | O_NONBLOCK | O_CLOEXEC | O_DIRECT | O_ACCMODE)
14
15 int SysOpen(size_t filename, size_t flags, size_t mode, size_t fdout, size_t) {
16     char path[400];
17     int fd;
18     int res;
19
20     if (flags & ~ALLOWABLE_FLAGS) {
21         return EINVAL;
22     }
23
24     if ((res = ReadStringFromUsermode(path, (const char*) filename, 399)) {
25         return res;
26     }
27
28     struct open_file* file;
29     if ((res = OpenFile(path, flags, mode, &file)) {
30         return res;
31     }
32
33     struct filedes_table* table = GetFileDescriptorTable(GetProcess());
34     if ((res = CreateFileDescriptor(table, file, &fd, flags & O_CLOEXEC)) {
35         CloseFile(file);
36         return res;
37     }
38
39     if ((res = WriteWordToUsermode((size_t*) fdout, fd)) {
40         CloseFile(file);
41         RemoveFileDescriptor(table, file);
42         return res;
43     }
44
45     return 0;
46 }

```

File: ./sys/calls/close.c

```

1 #include <syscall.h>
2 #include <errno.h>
3 #include <_syscallnum.h>
4 #include <thread.h>
5 #include <log.h>
6 #include <vfs.h>
7 #include <process.h>
8 #include <filedes.h>
9
10 int SysClose(size_t fd, size_t, size_t, size_t, size_t) {
11     struct open_file* file;
12     int res;
13
14     if ((res = GetFileFromDescriptor(GetFileDescriptorTable(GetProcess()), fd, &file)) {
15         return res;
16     }
17
18     return CloseFile(file);
19 }

```

File: ./sys/DS_Store

[binary]

File: ./sys/syscalls.c

```
1 #include <syscall.h>
2 #include <log.h>
3 #include <errno.h>
4 #include <_syscallnum.h>
5
6 typedef int (*system_call_t)(size_t, size_t, size_t, size_t, size_t);
7
8 const char* syscall_names[_SYSCALL_NUM_ENTRIES] = {
9     [SYSCALL_YIELD] = "yield",
10    [SYSCALL_TERMINATE] = "terminate",
11    [SYSCALL_MAPVIRT] = "map_virt",
12    [SYSCALL_UNMAPVIRT] = "unmap_virt",
13    [SYSCALL_OPEN] = "open",
14    [SYSCALL_READWRITE] = "read/write",
15    [SYSCALL_CLOSE] = "close",
16    [SYSCALL_SEEK] = "seek",
17    [SYSCALL_DUP] = "dup",
18    [SYSCALL_EXIT] = "exit",
19 };
20
21 static const system_call_t system_call_table[_SYSCALL_NUM_ENTRIES] = {
22     [SYSCALL_YIELD] = SysYield,
23     [SYSCALL_TERMINATE] = SysTerminate,
24     [SYSCALL_MAPVIRT] = SysMapVirt,
25     [SYSCALL_UNMAPVIRT] = SysUnmapVirt,
26     [SYSCALL_OPEN] = SysOpen,
27     [SYSCALL_READWRITE] = SysReadWrite,
28     [SYSCALL_CLOSE] = SysClose,
29     [SYSCALL_SEEK] = SysSeek,
30     [SYSCALL_DUP] = SysDup,
31     [SYSCALL_EXIT] = SysExit,
32 };
33
34 int HandleSystemCall(int call, size_t a, size_t b, size_t c, size_t d, size_t e) {
35     if (call >= _SYSCALL_NUM_ENTRIES) {
36         return ENOSYS;
37     }
38
39     return system_call_table[call](a, b, c, d, e);
40 }
```

File: ./arch/DS_Store

[binary]

File: ./arch/x86/DS_Store

[binary]

File: ./arch/x86/boot/kernel_entry.s

```
1
2 ;
3 ;
4 ; x86/kernel_entry.asm - Kernel Entry Point
5 ;
6 ; We want the kernel to conform to the Multiboot 2 standard. by doing this
7 ; the kernel can be loaded by common bootloaders, such as GRUB2, or directly
8 ; by the QEMU emulator
9 ;
10 ;
11 ;
12 ; Therefore, we must first define a few constants and flags required by Multiboot
13 ;
14 MBALIGN    equ 1 << 0
15 MEMINFO    equ 1 << 1
16 FLAGS      equ MBALIGN | MEMINFO
17 MAGIC      equ 0x1BADB002
18 CHECKSUM    equ ~(MAGIC + FLAGS)
19
20 ; We put the multiboot in a special section which gets placed at the start of
21 ; the binary. This allows the bootloader to find the Multiboot header.
22 ;
23 section .multiboot.data
24 align 4
25 dd MAGIC
26 dd FLAGS
27 dd CHECKSUM
28
29
30 ; We need a stack, so for the bootstrap process. We can define a quick 16KB
31 ; stack in the BSS section, which is always initialised to zeros
32 ;
33 section .bss
34 align 16
35 stack_bottom:
36 resb 4 * 1024
37 stack_top:
38
39
40 ; The kernel is being loaded to 0xC0100000. We need temporary bootstrap paging
41 ; structures handled here so that we can get the kernel to 0xC0100000 in virtual
42 ; memory.
43 ;
44 ; We will allocate one page directory, and one page table. With this, we can map
45 ; 4MB of memory. As the kernel starts at 1MB, we can actually have a kernel
46 ; of at most 3MB. We will replace these paging structures later once we get into
47 ; the proper kernel (and we'll release the physical memory behind it too!)
48 ;
49 align 4096
50 global boot_page_directory
51 global boot_page_table1
52 boot_page_directory: resb 4096
53 boot_page_table1: resb 4096
54
55 ; The start of the kernel itself - this will be called by the bootloader
56 ; We must place it in a special section so it appears at the start of the binary
57 ;
58 section .multiboot.text
59 global _start: function (_start, end - _start) ; calculate size of the _start function
60 extern KernelMain
61 extern _kernel_end
62 _start:
```

```

63 cli
64 cld
65
66 ; GRUB puts pointer in ebx, so we need to save it
67
68 ; Work out how many pages in the first 4MB need to be mapped
69 ; (we map the low 1MB, and then the kernel)
70 mov ecx, kernel_end
71 add ecx, 0xFFF
72 and ecx, 0xFFFFF000
73 shr ecx, 12
74 mov eax, 1024
75 sub eax, ecx
76
77 ; Get ready to loop over the page table
78 mov edi, boot_page_table1 - 0xC0000000
79 xor esi, esi
80 mov ecx, 1024 ; 1024 assumes 4MB of memory exists - we will only set the first 2MB as present
81 ; (as the kernel loads at 1MB, the kernel can be at most 1MB large)
82 ; (the page swapper will *hate you* if you 'invent' physical memory here)
83
84 .mapNextPage:
85 ; Combine the address with the present and writable flags
86 mov edx, esi
87 or edx, 3
88 cmp ecx, eax
89 jg .keep ; ** remember, the loop counter is going down **
90 xor edx, edx ; not present
91 .keep:
92 mov [edi], edx
93
94 .incrementPage:
95 ; Move onto the next page table entry, and the next corresponding physical page
96 add esi, 4096
97 add edi, 4
98 loop .mapNextPage
99
100 .endMapping:
101 ; Identity map and put the mappings at 0xC0000000
102 ; This way we won't page fault before we jump over to the kernel in high memory
103 ; (we are still in low memory)
104 mov [boot_page_directory - 0xC0000000 + 0], dword boot_page_table1 - 0xC0000000 + 3 + 256
105 mov [boot_page_directory - 0xC0000000 + 768 * 4], dword boot_page_table1 - 0xC0000000 + 3 + 256
106
107 ; Set the page directory
108 mov ecx, boot_page_directory - 0xC0000000
109 mov cr3, ecx
110
111 ; Enable paging
112 mov ecx, cr0
113 or ecx, (1 << 31)
114 or ecx, (1 << 16) ; enforce read-only pages in ring 0
115 mov cr0, ecx
116
117 ; This is why identity paging was required earlier, as paging is on, but we
118 ; are still in low memory (i.e. at 0x100000-ish)
119
120 ; Now jump to the higher half
121 lea ecx, KernelEntryPoint
122 jmp ecx
123 .end:
124
125 section .text
126
127 global vesa_pitch
128 global vesa_width
129 global vesa_height
130 global vesa_depth
131 global vesa_framebuffer
132
133 vesa_depth db 0
134 vesa_framebuffer dd 0
135 vesa_width dw 0
136 vesa_height dw 0
137 vesa_pitch dw 0
138
139 global x86_grub_table
140 x86_grub_table dd 0
141
142 ; The proper entry point of the kernel. Assumes the kernel is mapped into memory
143 ; at 0xC0100000.
144 KernelEntryPoint:
145 ; GRUB puts the address of a table in EBX, which we must use to find the
146 ; memory table. Note that we haven't trashed EBX up until this point.
147
148 ; TODO: kernel assumes the table is below 4MB, and that it is paged in
149 ; (which atm is only the case when it is below 1MB).
150 mov [x86_grub_table], ebx
151
152 ; Grab the video data the bootloader put into memory.
153 mov ax, [0x1000 + 16]
154 mov [vesa_pitch], ax
155
156 mov ax, [0x1000 + 18]
157 mov [vesa_width], ax
158
159 mov ax, [0x1000 + 20]
160 mov [vesa_height], ax
161
162 mov al, [0x1000 + 25]
163 mov [vesa_depth], al
164
165 mov eax, [0x1000 + 40]
166 mov [vesa_framebuffer], eax
167
168 ; Remove the identity paging and flush the TLB so the changes take effect
169 mov [boot_page_directory], dword 0
170 mov ecx, cr3
171 mov cr3, ecx
172
173 ; On x86, we'll store the current CPU number in the DR3 register (so user code cannot modify it)
174 ; Set it correctly now.
175 xor eax, eax
176 mov dr3, eax
177
178 ; Set the stack to the one we defined
179 mov esp, stack_top
180
181 ; Jump to the kernel main function
182 call KernelMain
183
184 ; We should never get here, but halt just in case
185 cli
186 hlt
187 jmp $

```

File: ./arch/x86/cpu/interrupt.c

```
1
2 #include <machine/regs.h>
3 #include <machine/interrupt.h>
4 #include <machine/pic.h>
5 #include <log.h>
6 #include <irq.h>
7 #include <irq1.h>
8 #include <virtual.h>
9 #include <syscall.h>
10 #include <panic.h>
11 #include <console.h>
12
13 #define ISR_SYSTEM_CALL 96
14 #define ISR_PAGE_FAULT 14
15 #define ISR_NMI 2
16
17 static bool ready_for_irqs = false;
18
19 static int GetRequiredIrql(int irq_num) {
20     if (irq_num == PIC_IRQ_BASE + 0) {
21         return IRQL_TIMER;
22     } else {
23         return IRQL_DRIVER + irq_num - PIC_IRQ_BASE;
24     }
25 }
26
27 void x86HandleInterrupt(struct x86_regs* r) {
28     int num = r->int_no;
29
30     if (num >= PIC_IRQ_BASE && num < PIC_IRQ_BASE + 16) {
31         RespondToIrq(num, GetRequiredIrql(num), r);
32     }
33     else if (num == ISR_PAGE_FAULT) {
34         extern size_t x86GetCr2();
35
36         int type = 0;
37         if (r->err_code & 1) {
38             type |= VM_READ;
39         }
40         if (r->err_code & 2) {
41             type |= VM_WRITE;
42         }
43         if (r->err_code & 4) {
44             type |= VM_USER;
45         }
46         if (r->err_code & 16) {
47             type |= VM_EXEC;
48         }
49
50         LogWriteSerial("\n\nPage fault: cr2 0x%X, eip 0x%X, nos-err 0x%X\n", x86GetCr2(), r->eip, type);
51
52         HandleVirtFault(x86GetCr2(), type);
53     }
54     else if (num == ISR_NMI) {
55         Panic(PANIC_NON_MASKABLE_INTERRUPT);
56     }
57     else if (num == ISR_SYSTEM_CALL) {
58         r->eax = HandleSystemCall(r->eax, r->ebx, r->ecx, r->edx, r->esi, r->edi);
59     }
60     else {
61         LogWriteSerial("Got interrupt %d. (r->eip = 0x%X)\n", num, r->eip);
62         UnhandledFault();
63     }
64 }
65
66 void ArchSendEoi(int irq_num) {
67     SendPicEoi(irq_num);
68 }
69
70 void ArchSetIrql(int irql) {
71     if (irql == IRQL_HIGH || irql == IRQL_TIMER || !x86IsReadyForIrqs()) {
72         /*
73          * Interrupts stay off.
74          */
75         return;
76     }
77
78     if (irql >= IRQL_DRIVER) {
79         int irq_num = irql - IRQL_DRIVER;
80
81         /*
82          * We want to disable all higher IRQs (as the PIC puts the lowest priority interrupts at
83          * high numbers), as well as our self. Allow IRQ2 to stay enabled as it is used internally.
84          */
85         uint16_t mask = (0xFFFF ^ ((1 << irq_num) - 1)) & ~(1 << 2);
86         DisablePicLines(mask);
87     }
88     else {
89         /*
90          * Allow everything to go through.
91          */
92         DisablePicLines(0x0000);
93     }
94
95     ArchEnableInterrupts();
96 }
97
98 bool x86IsReadyForIrqs(void) {
99     return ready_for_irqs;
100 }
101
102 void x86MakeReadyForIrqs(void) {
103     ready_for_irqs = true;
104     RaiseIrql(GetIrql());
105 }
```

File: ./arch/x86/cpu/cpu.c

```

1
2 #include <stdbool.h>
3 #include <virtual.h>
4 #include <machine/gdt.h>
5 #include <machine/idt.h>
6 #include <machine/tss.h>
7 #include <machine/pic.h>
8 #include <machine/pit.h>
9 #include <cpu.h>
10 #include <machine/portio.h>
11 #include <machine/interrupt.h>
12 #include <errno.h>
13 #include <driver.h>
14
15 static void x86EnableNMIs(void) {
16     outb(0x70, inb(0x70) & 0x7F);
17     inb(0x71);
18 }
19
20 void ArchInitBootstrapCpu(struct cpu*) {
21     x86InitGdt();
22     x86InitIdt();
23     x86InitTss();
24
25     InitPic();
26     InitPit(40);
27
28     ArchEnableInterrupts();
29     x86MakeReadyForIrqs();
30     x86EnableNMIs();
31 }
32
33 bool ArchInitNextCpu(struct cpu*) {
34     return false;
35 }
36
37 static void x86Reboot(void) {
38     uint8_t good = 0x02;
39     while (good < 0x02) {
40         good = inb(0x64);
41     }
42     outb(0x64, 0xFE);
43 }
44
45 static void x86Shutdown(void) {
46     size_t acpicaShutdown = GetSymbolAddress("AcpicaShutdown");
47     if (acpicaShutdown != 0) {
48         ((void (*)(void)) acpicaShutdown)();
49     }
50
51     /*
52      * Some emulators have ways of doing a shutdown if we don't have ACPI support yet.
53      */
54     outw(0xB004, 0x2000); // Bochs and old QEMU
55     outw(0x0604, 0x2000); // New QEMU
56     outw(0x4004, 0x3400); // VirtualBox
57     outw(0x0600, 0x0034); // Cloud Hypervisor
58 }
59
60 static void x86Sleep(void) {
61     size_t acpicaSleep = GetSymbolAddress("AcpicaSleep");
62     if (acpicaSleep != 0) {
63         ((void (*)(void)) acpicaSleep)();
64     }
65 }
66
67 int ArchSetPowerState(int power_state) {
68     switch (power_state) {
69     case ARCH_POWER_STATE_REBOOT:
70         x86Reboot();
71         break;
72     case ARCH_POWER_STATE_SHUTDOWN:
73         x86Shutdown();
74         break;
75     case ARCH_POWER_STATE_SLEEP: {
76         x86Sleep();
77         break;
78     }
79     default:
80         return EINVAL;
81     }
82
83     while (1) {
84         ArchStallProcessor();
85     }
86 }

```

File: /arch/x86/include/idt.h

```
#pragma once
```

```
#include <common.h>
```

```
/* x86/lowlevel/idt.h - Interrupt Descriptor Table
```

```
 *
 *
 */
```

```
/*
```

```
 * An entry in the IDT. The offset is the address the CPU will jump to,
 * and the selector is what segment should be used (i.e. we need to have
 * setup a GDT already). The layout of this structure is mandated by the CPU.
 */
```

```
struct idt_entry
{
    uint16_t isr_offset_low;
    uint16_t segment_selector;
    uint8_t reserved;
    uint8_t type;
    uint16_t isr_offset_high;
} __attribute__((packed));
```

```
/*
```

```
 * Used to tell the CPU where the IDT is and how long it is.
 * The layout of this structure is mandated by the CPU.
 */
```

```
struct idt_ptr
{
    uint16_t size;
```

```
size_t location;
} __attribute__((packed));
```

```
void x86InitIdt(void);
```

File: */arch/x86/include/config.h*

```
#pragma once

/*
 * As this is for x86 (not x86-64), we set the limit to 4GB. On x86-64, we can set
 * it larger. This will make the bitmap much larger, but this is no problem on an
 * x86-64 system (only ancient x86 systems will have e.g. 4MB of RAM).
 */
#define ARCH_MAX_RAM_KBS (1024 * 1024 * 4)

#define ARCH_PAGE_SIZE 4096

/*
 * Non-inclusive of ARCH_USER_AREA_LIMIT
 */
#define ARCH_USER_AREA_BASE 0x08000000
#define ARCH_USER_AREA_LIMIT 0xC0000000

#define ARCH_USER_STACK_BASE 0x08000000
#define ARCH_USER_STACK_LIMIT 0x10000000

/*
 * Non-inclusive of ARCH_KRNL_SBRK_LIMIT. Note that we can't use the top 8MB,
 * as we use that for recursive mapping.
 */
#define ARCH_KRNL_SBRK_BASE 0xC4000000
#define ARCH_KRNL_SBRK_LIMIT 0xFFC00000
#define ARCH_PROG_LOADER_BASE 0xBFC00000
#define ARCH_PROG_LOADER_ENTRY 0xBFC00000

#define ARCH_MAX_CPU_ALLOWED 16

#undef ARCH_BIG_ENDIAN
#define ARCH_LITTLE_ENDIAN

#include <machine/gdt.h>
#include <machine/idt.h>
#include <machine/tss.h>
#include <machine/regs.h>

typedef struct {
    /* Plz keep tss at the top, thread switching assembly needs it */
    struct tss* tss;

    struct gdt_entry gdt[16];
    struct idt_entry idt[256];

    struct gdt_ptr gdt_r;
    struct idt_ptr idt_r;

} platform_cpu_data_t;

typedef struct {
    size_t p_page_directory; // cr3
    size_t* v_page_directory; // what we use to access the tables

} platform_vas_data_t;

typedef struct x86_regs platform_irq_context_t;
```

File: */arch/x86/include/virtual.h*

```
#pragma once

#include <stddef.h>

__attribute__((fastcall)) size_t x86KernelMemoryToPhysical(size_t virtual);
```

File: */arch/x86/include/pit.h*

```
#pragma once

void InitPit(int hertz);
```

File: */arch/x86/include/dev.h*

```
#pragma once

void InitIde(void);
void InitFloppy(void);
```

File: */arch/x86/include/elf.h*

```
#pragma once
#include <stdint.h>
#include <stddef.h>
#include <stdbool.h>

#define ELF_NIDENT 16

typedef uint16_t Elf32_Half; // Unsigned half int
typedef uint32_t Elf32_Off; // Unsigned offset
typedef uint32_t Elf32_Addr; // Unsigned address
```

```

typedef uint32_t Elf32_Word; // Unsigned int
typedef int32_t Elf32_Sword; // Signed int

struct Elf32_Ehdr
{
    uint8_t e_ident[ELF_NIDENT];
    Elf32_Half e_type;
    Elf32_Half e_machine;
    Elf32_Word e_version;
    Elf32_Addr e_entry;
    Elf32_Off e_phoff;
    Elf32_Off e_shoff;
    Elf32_Word e_flags;
    Elf32_Half e_ehsize;
    Elf32_Half e_phsize;
    Elf32_Half e_phnum;
    Elf32_Half e_shsize;
    Elf32_Half e_shnum;
    Elf32_Half e_shstrndx;
};

enum Elf_Ident
{
    EI_MAG0 = 0, // 0x7F
    EI_MAG1 = 1, // 'E'
    EI_MAG2 = 2, // 'L'
    EI_MAG3 = 3, // 'F'
    EI_CLASS = 4, // Architecture (32/64)
    EI_DATA = 5, // Byte Order
    EI_VERSION = 6, // ELF Version
    EI_OSABI = 7, // OS Specific
    EI_ABIVERSION = 8, // OS Specific
    EI_PAD = 9 // Padding
};

#define ELFMAG0 0x7F // e_ident[EI_MAG0]
#define ELFMAG1 'E' // e_ident[EI_MAG1]
#define ELFMAG2 'L' // e_ident[EI_MAG2]
#define ELFMAG3 'F' // e_ident[EI_MAG3]

#define ELFDATA2LSB (1) // Little Endian
#define ELFCLASS32 (1) // 32-bit Architecture

enum Elf_Type
{
    ET_NONE = 0, // Unkown Type
    ET_REL = 1, // Relocatable File
    ET_EXEC = 2 // Executable File
};

#define EM_386 (3) // x86 Machine Type
#define EV_CURRENT (1) // ELF Current Version

struct Elf32_Shdr
{
    Elf32_Word sh_name;
    Elf32_Word sh_type;
    Elf32_Word sh_flags;
    Elf32_Addr sh_addr;
    Elf32_Off sh_offset;
    Elf32_Word sh_size;
    Elf32_Word sh_link;
    Elf32_Word sh_info;
    Elf32_Word sh_addralign;
    Elf32_Word sh_entsize;
};

#define SHN_UNDEF 0x0000 // Undefined/Not present
#define SHN_ABS 0xFFFF1 // Absolute value

enum ShT_Types
{
    SHT_NULL = 0, // Null section
    SHT_PROGBITS = 1, // Program information
    SHT_SYMTAB = 2, // Symbol table
    SHT_STRTAB = 3, // String table
    SHT_RELA = 4, // Relocation (w/ addend)
    SHT_NOBITS = 8, // Not present in file
    SHT_REL = 9, // Relocation (no addend)
};

enum ShT_Attributes
{
    SHF_WRITE = 0x01, // Writable section
    SHF_ALLOC = 0x02 // Exists in memory
};

struct Elf32_Sym
{
    Elf32_Word st_name;
    Elf32_Addr st_value;
    Elf32_Word st_size;
    uint8_t st_info;
    uint8_t st_other;
    Elf32_Half st_shndx;
};

#define ELF32_ST_BIND(INFO) ((INFO) >> 4)
#define ELF32_ST_TYPE(INFO) ((INFO) & 0x0F)

enum StT_Bindings
{
    STB_LOCAL = 0, // Local scope
    STB_GLOBAL = 1, // Global scope
    STB_WEAK = 2 // Weak, (ie. __attribute__((weak)))
};

```



```

enum StT_Types
{
    STT_NOTYPE = 0, // No type
    STT_OBJECT = 1, // Variables, arrays, etc.
    STT_FUNC = 2 // Methods or functions
};

struct Elf32_Rel
{
    Elf32_Addr r_offset;
    Elf32_Word r_info;
};

struct Elf32_Rela {
    Elf32_Addr r_offset;
    Elf32_Word r_info;
    Elf32_Sword r_addend;
};

#define ELF32_R_SYM(INFO) ((INFO) >> 8)
#define ELF32_R_TYPE(INFO) ((uint8_t)(INFO))

enum RtT_Types
{
    R_386_NONE = 0, // No relocation
    R_386_32 = 1, // Symbol + Offset
    R_386_PC32 = 2, // Symbol + Offset - Section Offset
    R_386_RELATIVE = 8,
};

struct Elf32_Phdr
{
    Elf32_Word p_type;
    Elf32_Off p_offset;
    Elf32_Addr p_vaddr;
    Elf32_Addr p_paddr;
    Elf32_Word p_filesz;
    Elf32_Word p_memsz;
    Elf32_Word p_flags;
    Elf32_Word p_align;
};

enum PH_Types
{
    PHT_NULL = 0,
    PHT_LOAD = 1,
    PHT_DYNAMIC = 2,
    PHT_INTERP = 3,
    PHT_NOTE = 4,
    PHT_SHLIB = 5,
    PHT_PHDR = 6,
    PHT_LOOS = 0x60000000,
    PHT_HIOS = 0x6FFFFFFF,
    PHT_LOPROC = 0x70000000,
    PHT_HIPROC = 0x7FFFFFFF
};

#define ELF32_R_SYM(INFO) ((INFO) >> 8)
#define ELF32_R_TYPE(INFO) ((uint8_t)(INFO))

#define DO_386_32(S, A) ((S) + (A))
#define DO_386_RELATIVE(B, A) ((B) + (A))
#define DO_386_PC32(S, A, P) ((S) + (A) - (P))

#define PF_X 1
#define PF_W 2
#define PF_R 4

```

File: */arch/x86/include/pic.h*

```

#pragma once

#include <stdbool.h>
#include <stdint.h>

#define PIC_IRQ_BASE 32

void InitPic(void);
void SendPicEoi(int irq_num);
bool IsPicIrqSpurious(int irq_num);
void DisablePicLines(uint16_t irq_bitfield);

```

File: */arch/x86/include/interrupt.h*

```

#pragma once

#include <stdbool.h>

bool x86IsReadyForIrqs(void);
void x86MakeReadyForIrqs(void);

```

File: */arch/x86/include/regs.h*

```

#pragma once
#include <common.h>

struct x86_regs
{
    /*
     * The registers that are pushed to us in x86/lowlevel/trap.s
     *
     * SS is the first thing pushed, and thus the last to be popped
     * GS is the last thing pushed, and thus the first to be popped
     */
};

```

```

*/
size_t gs, fs, es, ds;
size_t edi, esi, ebp, esp, ebx, edx, ecx, eax;
size_t int_no, err_code;
size_t eip, cs, eflags, useresp, ss;
};

```

File: `/arch/x86/include/tss.h`

```

#pragma once

/* x86/lowlevel/tss.h - Task State Segment
 *
 *
 */

#include <common.h>

/*
 * The task state segment was designed to store information about a task so
 * that task switching could be done in hardware. We do not use it for this purpose,
 * instead only using it to set the stack correctly after a user -> kernel switch.
 *
 * The layout of this structure is mandated by the CPU.
 */
struct tss
{
    uint16_t link;
    uint16_t reserved0;
    uint32_t esp0;
    uint16_t ss0;
    uint16_t reserved1;
    uint32_t esp1;
    uint16_t ss1;
    uint16_t reserved2;
    uint32_t esp2;
    uint16_t ss2;
    uint16_t reserved3;
    uint32_t cr3;
    uint32_t eip;
    uint32_t eflags;
    uint32_t eax;
    uint32_t ecx;
    uint32_t edx;
    uint32_t ebx;
    uint32_t esp;
    uint32_t ebp;
    uint32_t esi;
    uint32_t edi;
    uint16_t es;
    uint16_t reserved4;
    uint16_t cs;
    uint16_t reserved5;
    uint16_t ss;
    uint16_t reserved6;
    uint16_t ds;
    uint16_t reserved7;
    uint16_t fs;
    uint16_t reserved8;
    uint16_t gs;
    uint16_t reserved9;
    uint16_t ldtr;
    uint16_t reserved10;
    uint16_t reserved11;
    uint16_t ioph;
} __attribute__((packed));

void x86InitTss(void);

```

File: `/arch/x86/include/gdt.h`

```

#pragma once

#include <common.h>

/*
 * An entry in the GDT table. The layout of this structure is mandated by the CPU.
 */
struct gdt_entry
{
    uint16_t limit_low;
    uint16_t base_low;
    uint8_t base_middle;
    uint8_t access;
    uint8_t flags_and_limit_high;
    uint8_t base_high;
} __attribute__((packed));

/*
 * Describes the GDT address and size. We use the address of this structure
 * to tell the CPU where the GDT is. The layout of this structure is mandated by the CPU.
 */
struct gdt_ptr
{
    uint16_t size;
    size_t location;
} __attribute__((packed));

struct tss;

void x86InitGdt(void);
uint16_t x86AddTssToGdt(struct tss* tss);

```

File: /arch/x86/include/portio.h

```
#pragma once

/*
 * x86/portio.h - Port Input / Output
 *
 *
 * On the x86, a lot of older hardware is accessed using IO ports. A port has an
 * address from 0x0000 to 0xFFFF, and can be read from and written to using special
 * instructions.
 *
 * We are going to inline these functions, as they are all single instructions.
 */

#include <common.h>
#include <log.h>

/*
 * Writing to ports
 */
always_inline void outb(uint16_t port, uint8_t value)
{
    asm volatile ("outb %0, %1" : "a"(value), "Nd"(port));
}

always_inline void outw(uint16_t port, uint16_t value)
{
    asm volatile ("outw %0, %1" : "a"(value), "Nd"(port));
}

always_inline void outl(uint16_t port, uint32_t value)
{
    asm volatile ("outl %0, %1" : "a"(value), "Nd"(port));
}

/*
 * Reading from ports
 */
always_inline uint8_t inb(uint16_t port)
{
    uint8_t value;
    asm volatile ("inb %1, %0"
        : "=a"(value)
        : "Nd"(port));
    return value;
}

always_inline uint16_t inw(uint16_t port)
{
    uint16_t value;
    asm volatile ("inw %1, %0"
        : "=a"(value)
        : "Nd"(port));
    return value;
}

always_inline uint32_t inl(uint16_t port)
{
    uint32_t value;
    asm volatile ("inl %1, %0"
        : "=a"(value)
        : "Nd"(port));
    return value;
}
```

File: /arch/x86/driver.ld

```
OUTPUT_FORMAT("elf32-i386")

SECTIONS
{
    . = 0xD0000000;

    .text BLOCK(4096) : ALIGN(4096)
    {
        *(.text)
        *(.rodata)
    }

    .data BLOCK(4096) : ALIGN(4096)
    {
        *(.data)
    }

    .lockedtext BLOCK(4096) : ALIGN(4096)
    {
        *(.lockedtext)
        *(.lockedrodata)
    }

    .lockeddata BLOCK(4096) : ALIGN(4096)
    {
        *(.lockeddata)
    }

    .bss BLOCK(4096) : ALIGN(4096)
    {
        *(COMMON)
        *(.bss)
        *(.bootstrap_stack)
    }
}
```

```
/DISCARD/ :  
{  
  *(.comment)  
}
```

File: `/arch/x86/dev/pic.c`

```

1  #include <machine/pic.h>
2  #include <machine/portio.h>
3  #include <arch.h>
4
5  #define PIC1_COMMAND    0x20
6  #define PIC1_DATA       0x21
7  #define PIC2_COMMAND    0xA0
8  #define PIC2_DATA       0xA1
9
10 #define PIC_EOI         0x20
11 #define PIC_REG_ISR     0x0B
12
13 #define ICW1_ICW4       0x01
14 #define ICW1_INIT       0x10
15 #define ICW4_8086       0x01
16
17 /*
18  * Delay for a short period of time, for use in between IO calls to the PIC.
19  * This is required as some PICs have a hard time keeping up with the speed
20  * of modern CPUs (the original PIC was introduced in 1976!).
21  */
22 static void IoWait(void) {
23     asm volatile ("nop");
24 }
25
26 /*
27  * Read an internal PIC register.
28  */
29 static uint16_t ReadPicReg(int ocw3) {
30     outb(PIC1_COMMAND, ocw3);
31     outb(PIC2_COMMAND, ocw3);
32     return ((uint16_t) inb(PIC2_COMMAND) << 8) | inb(PIC1_COMMAND);
33 }
34
35 /*
36  * Due to a race condition between the PIC and the CPU, we sometimes get a
37  * 'spurious' interrupt sent to the CPU on IRQ 7 or 15. Distinguishing them is
38  * important - we don't need to send an EOI after a spurious interrupt.
39  */
40 bool IsPicIrqSpurious(int irq_num) {
41     if (irq_num == PIC_IRQ_BASE + 7) {
42         uint16_t isr = ReadPicReg(PIC_REG_ISR);
43         return !(isr & (1 << 7));
44     } else if (irq_num == PIC_IRQ_BASE + 15) {
45         uint16_t isr = ReadPicReg(PIC_REG_ISR);
46         if (!(isr & (1 << 15))) {
47             /*
48              * It is spurious, but the primary PIC doesn't know that, as it came
49              * from the secondary PIC. So only send an EOI to the primary PIC.
50              */
51             outb(PIC1_COMMAND, PIC_EOI);
52             return true;
53         }
54     }
55     return false;
56 }
57
58 /*
59  * Acknowledge the previous interrupt. We will not receive any interrupts of
60  * the same type until we have acknowledged it.
61  */
62 void SendPicEoi(int irq_num) {
63     if (irq_num >= PIC_IRQ_BASE + 8) {
64         outb(PIC2_COMMAND, PIC_EOI);
65     }
66     outb(PIC1_COMMAND, PIC_EOI);
67 }
68
69 /*
70  * Change which interrupt numbers are used by the IRQs. They will initially
71  * use interrupts 0 through 15, which isn't very good as it conflicts with the
72  * interrupt numbers for the CPU exceptions.
73  */
74 static void RemapPic(int offset) {
75     uint8_t mask1 = inb(PIC1_DATA);
76     uint8_t mask2 = inb(PIC2_DATA);
77
78     outb(PIC1_COMMAND, ICW1_INIT | ICW1_ICW4);
79     IoWait();
80     outb(PIC2_COMMAND, ICW1_INIT | ICW1_ICW4);
81     IoWait();
82     outb(PIC1_DATA, offset);
83     IoWait();
84     outb(PIC2_DATA, offset + 8);
85     IoWait();
86     outb(PIC1_DATA, 4);
87     IoWait();
88     outb(PIC2_DATA, 2);
89     IoWait();
90
91     outb(PIC1_DATA, ICW4_8086);
92     IoWait();
93     outb(PIC2_DATA, ICW4_8086);
94     IoWait();
95
96     outb(PIC1_DATA, mask1);
97     outb(PIC2_DATA, mask2);
98 }
99
100 /**
101  * Set which IRQ numbers are disabled. Overwrites the previous call to this
102  * function completely (i.e. this is 'equals', not an 'and' or 'or'.) To disable
103  * all lines, specify 0xFFFF. To enable all lines, specify 0x0000.
104  */
105 void DisablePicLines(uint16_t irq_bitfield) {
106     static uint16_t prev = 0xFFFF;
107     if (prev != irq_bitfield) {
108         outb(PIC1_DATA, irq_bitfield & 0xFF);
109         outb(PIC2_DATA, irq_bitfield >> 8);
110         prev = irq_bitfield;
111     }
112 }
113
114 void InitPic(void) {
115     RemapPic(PIC_IRQ_BASE);
116     DisablePicLines(0x0000);
117 }

```

File: ./arch/x86/dev/init.c

```

1 #include <machine/dev.h>
2 #include <machine/portio.h>
3 #include <driver.h>
4 #include <thread.h>
5 #include <panic.h>
6 #include <log.h>
7 #include <virtual.h>
8
9 static size_t Loadx86Driver(const char* filename, const char* init) {
10     if (RequireDriver(filename)) {
11         PanicEx(PANIC_REQUIRED_DRIVER_NOT_FOUND, filename);
12     }
13     if (GetSymbolAddress(init) != 0) {
14         PanicEx(PANIC_REQUIRED_DRIVER_MISSING_SYMBOL, filename);
15     }
16     return addr;
17 }
18
19 static void LoadSlowDriversInBackground(void*) {
20     ((void(*)()) (Loadx86Driver("sys:/acpi.sys", "InitAcpi")))();
21 }
22
23 void ArchInitDev(bool fs) {
24     if (!fs)
25         InitIde();
26     //InitFloppy();
27
28     else {
29         ((void(*)()) (Loadx86Driver("sys:/vesa.sys", "InitVesa")))();
30         ((void(*)()) (Loadx86Driver("sys:/ps2.sys", "InitPs2")))();
31         CreateThread(LoadSlowDriversInBackground, NULL, GetVas(), "drvloader");
32     }
33 }

```

File: /arch/x86/dev/floppy.c

```

1 #include <common.h>
2 #include <semaphore.h>
3 #include <log.h>
4 #include <thread.h>
5 #include <vfs.h>
6 #include <string.h>
7 #include <transfer.h>
8 #include <assert.h>
9 #include <irq.h>
10 #include <errno.h>
11 #include <machine/pic.h>
12 #include <machine/portio.h>
13 #include <heap.h>
14 #include <virtual.h>
15 #include <stdlib.h>
16 #include <sys/stat.h>
17 #include <dirent.h>
18 #include <irq.h>
19 #include <diskutil.h>
20
21 #define CYLINDER_SIZE (512 * 18 * 2)
22
23 #define FLOPPY_DOR 2
24 #define FLOPPY_MSR 4
25 #define FLOPPY_FIFO 5
26 #define FLOPPY_CCR 7
27
28 #define CMD_SPECIFY 0x03
29 #define CMD_READ 0x06
30 #define CMD_RECALIBRATE 0x07
31 #define CMD_SENSE_INT 0x08
32 #define CMD_SEEK 0x0F
33 #define CMD_CONFIGURE 0x13
34
35 struct semaphore* floppy_lock = NULL;
36
37 struct floppy_data {
38     int disk_num;
39     uint8_t* cylinder_buffer;
40     uint8_t* cylinder_zero;
41     size_t base;
42     struct disk_partition_helper partitions;
43     int stored_cylinder;
44     bool got_cylinder_zero;
45 };
46
47 static void FloppyWriteCommand(struct floppy_data* fdp, int cmd) {
48     int base = fdp->base;
49
50     for (int i = 0; i < 60; ++i) {
51         SleepMilli(10);
52         if (inb(base + FLOPPY_MSR) & 0x80) {
53             outb(base + FLOPPY_FIFO, cmd);
54             return;
55         }
56     }
57 }
58
59 static uint8_t FloppyReadData(struct floppy_data* fdp) {
60     int base = fdp->base;
61     for (int i = 0; i < 60; ++i) {
62         SleepMilli(10);
63         if (inb(base + FLOPPY_MSR) & 0x80) {
64             return inb(base + FLOPPY_FIFO);
65         }
66     }
67
68     return 0;
69 }
70
71 static void FloppyCheckInterrupt(struct floppy_data* fdp, int* st0, int* cyl) {
72     FloppyWriteCommand(fdp, CMD_SENSE_INT);
73     *st0 = FloppyReadData(fdp);
74     *cyl = FloppyReadData(fdp);
75 }
76
77 /*
78  * The state can be 0 (off), 1 (on) or 2 (currently on, but will shortly be turned off).
79  */
80 static volatile int floppy_motor_state = 0;
81 static volatile int floppy_motor_ticks = 0;
82
83 static void FloppyMotor(struct floppy_data* fdp, bool state) {
84     if (state)
85         if (!floppy_motor_state) {
86             outb(fdp->base + FLOPPY_DOR, 0x1C);
87             SleepMilli(150);
88         }
89     floppy_motor_state = 1;
90 }

```

```

91     } else {
92         floppy_motor_state = 2;
93         floppy_motor_ticks = 1000;
94     }
95 }
96
97 static volatile bool floppy_got_irq = false;
98
99 static int FloppyIrqWait(void) {
100     int timeout = 0;
101     while (!floppy_got_irq) {
102         SleepMilli(10);
103         if (++timeout > 200) {
104             return ETIMEDOUT;
105         }
106     }
107
108     floppy_got_irq = false;
109     return 0;
110 }
111
112 static int FloppyIrqHandler(struct x86_regs*) {
113     floppy_got_irq = true;
114     return 0;
115 }
116
117 static void FloppyMotorControlThread(void*) {
118     while (1) {
119         SleepMilli(50);
120         if (floppy_motor_state == 2) {
121             floppy_motor_ticks -= 50;
122             if (floppy_motor_ticks <= 0) {
123                 /*
124                  * Actually turn off the motor.
125                  */
126                 outb(0x3F0 + FLOPPY_DOR, 0x0C);
127                 floppy_motor_state = 0;
128             }
129         }
130     }
131 }
132
133 /*
134 * Move to cylinder 0.
135 */
136 static int FloppyCalibrate(struct floppy_data* flp) {
137     LogWriteSerial("FloppyCalibrate\n");
138
139     int st0 = -1;
140     int cyl = -1;
141
142     FloppyMotor(flp, true);
143
144     for (int i = 0; i < 10; ++i) {
145         FloppyWriteCommand(flp, CMD_RECALIBRATE);
146         FloppyWriteCommand(flp, 0);
147
148         FloppyIrqWait();
149         FloppyCheckInterrupt(flp, &st0, &cyl);
150
151         if (st0 & 0xC0) {
152             continue;
153         }
154
155         if (cyl == 0) {
156             FloppyMotor(flp, false);
157             return 0;
158         }
159     }
160
161     FloppyMotor(flp, false);
162     return EIO;
163 }
164
165 static void FloppyConfigure(struct floppy_data* flp) {
166     FloppyWriteCommand(flp, CMD_CONFIGURE);
167     FloppyWriteCommand(flp, 0x00);
168     FloppyWriteCommand(flp, 0x08);
169     FloppyWriteCommand(flp, 0x00);
170 }
171
172 static int FloppyReset(struct floppy_data* flp) {
173     int base = flp->base;
174     outb(base + FLOPPY_DOR, 0x00);
175     outb(base + FLOPPY_DOR, 0x0C);
176
177     FloppyIrqWait();
178
179     for (int i = 0; i < 4; ++i) {
180         int st0, cyl;
181         FloppyCheckInterrupt(flp, &st0, &cyl);
182     }
183
184     outb(base + FLOPPY_CCR, 0x00);
185
186     FloppyMotor(flp, true);
187
188     FloppyWriteCommand(flp, CMD_SPECIFY);
189     FloppyWriteCommand(flp, 0x0F);
190     FloppyWriteCommand(flp, 0x02);
191
192     SleepMilli(300);
193     FloppyConfigure(flp);
194     SleepMilli(300);
195     FloppyMotor(flp, false);
196
197     return FloppyCalibrate(flp);
198 }
199
200 static int FloppySeek(struct floppy_data* flp, int cylinder, int head) {
201     int st0, cyl;
202     FloppyMotor(flp, true);
203
204     for (int i = 0; i < 10; ++i) {
205         FloppyWriteCommand(flp, CMD_SEEK);
206         FloppyWriteCommand(flp, head < 2);
207         FloppyWriteCommand(flp, cylinder);
208
209         FloppyIrqWait();
210         FloppyCheckInterrupt(flp, &st0, &cyl);
211
212         if (st0 & 0xC0) {
213             continue;
214         }
215
216         if (cyl == cylinder) {
217             FloppyMotor(flp, false);
218             return 0;
219         }
220     }

```

```

221
222     FloppyMotor(flp, false);
223     return EIO;
224 }
225
226 static void FloppyDmaInit(void) {
227     /*
228      * Put the data at *physical address* 0x10000. The address can be anywhere
229      * under 24MB that doesn't cross a 64KB boundary. We choose this location as
230      * it should be unused as this is where the temporary copy of the kernel was
231      * stored during boot.
232      */
233     uint32_t addr = (uint32_t) 0x10000;
234
235     /*
236      * We must give the DMA the actual count minus 1.
237      */
238     int count = 0x4800 - 1;
239
240     /*
241      * Send some magical stuff to the DMA controller.
242      */
243     outb(0x0A, 0x06);
244     outb(0x0C, 0xFF);
245     outb(0x04, (addr >> 0) & 0xFF);
246     outb(0x04, (addr >> 8) & 0xFF);
247     outb(0x81, (addr >> 16) & 0xFF);
248     outb(0x0C, 0xFF);
249     outb(0x05, (count >> 0) & 0xFF);
250     outb(0x05, (count >> 8) & 0xFF);
251     outb(0x0B, 0x46);
252     outb(0x0A, 0x02);
253 }
254
255 static int FloppyDoCylinder(struct floppy_data* flp, int cylinder) {
256     /*
257      * Move both heads to the correct cylinder.
258      */
259     if (FloppySeek(flp, cylinder, 0) != 0) return EIO;
260     if (FloppySeek(flp, cylinder, 1) != 0) return EIO;
261
262     /*
263      * This time, we'll try up to 20 times.
264      */
265     for (int i = 0; i < 20; ++i) {
266         FloppyMotor(flp, true);
267
268         if (i % 5 == 3) {
269             if (i % 10 == 8) {
270                 FloppyReset(flp);
271                 FloppyMotor(flp, true);
272             }
273
274             FloppyCalibrate(flp);
275
276             if (FloppySeek(flp, cylinder, 0) != 0) return EIO;
277             if (FloppySeek(flp, cylinder, 1) != 0) return EIO;
278         }
279
280         FloppyDmaInit();
281
282         SleepMilli(100);
283
284         /*
285          * Send the read command.
286          */
287         FloppyWriteCommand(flp, CMD_READ | 0xC0);
288         FloppyWriteCommand(flp, 0);
289         FloppyWriteCommand(flp, cylinder);
290         FloppyWriteCommand(flp, 0);
291         FloppyWriteCommand(flp, 1);
292         FloppyWriteCommand(flp, 2);
293         FloppyWriteCommand(flp, 18);
294         FloppyWriteCommand(flp, 0x1B);
295         FloppyWriteCommand(flp, 0xFF);
296
297         FloppyIrqWait();
298
299         /*
300          * Read back some status information, some of which is very mysterious.
301          */
302         uint8_t st0, st1, st2, rcy, rhe, rse, bps;
303         st0 = FloppyReadData(flp);
304         st1 = FloppyReadData(flp);
305         st2 = FloppyReadData(flp);
306         rcy = FloppyReadData(flp);
307         rhe = FloppyReadData(flp);
308         rse = FloppyReadData(flp);
309         bps = FloppyReadData(flp);
310
311         /*
312          * Check for errors. More tests can be done, but it would make the code
313          * even longer.
314          */
315         if (st0 & 0xC0) {
316             continue;
317         }
318         if (st1 & 0x80) {
319             continue;
320         }
321         if (st0 & 0x8) {
322             continue;
323         }
324         if (st1 & 0x20) {
325             continue;
326         }
327         if (st1 & 0x10) {
328             continue;
329         }
330         if (bps != 2) {
331             continue;
332         }
333
334         (void) st2;
335         (void) rcy;
336         (void) rhe;
337         (void) rse;
338
339         FloppyMotor(flp, false);
340         return 0;
341     }
342
343     FloppyMotor(flp, false);
344     return EIO;
345 }
346
347 static int FloppyIo(struct floppy_data* flp, struct transfer* io) {
348     EXACT_IRQL(IRQL_STANDARD);
349
350     if (io->direction == TRANSFER_WRITE) {

```



```

351     return EROFS;
352 }
353
354 int lba = io->offset / 512;
355 int count = io->length_remaining / 512;
356
357 if (io->offset % 512 != 0) {
358     return EINVAL;
359 }
360 if (io->length_remaining % 512 != 0) {
361     return EINVAL;
362 }
363 if (count <= 0 || count > 0xFF || lba < 0 || lba >= 2880) {
364     return EINVAL;
365 }
366
367 AcquireMutex(floppy_lock, -1);
368
369 next_sector:;
370 /*
371  * Floppies use CHS (cylinder, head, sector) for addressing sectors instead
372  * of LBA (linear block addressing). Hence we need to convert to CHS.
373  * Note that sector is 1-based, whereas cylinder and head are 0-based.
374  */
375 int head = (lba % (18 * 2)) / 18;
376 int cylinder = (lba / (18 * 2));
377 int sector = (lba % 18) + 1;
378
379 /*
380  * Cylinder 0 has some commonly used data, so cache it separately for
381  * improved speed.
382  */
383 if (cylinder == 0) {
384     if (!flp->got_cylinder_zero) {
385         flp->got_cylinder_zero = true;
386
387         int error = FloppyDoCylinder(flp, cylinder);
388
389         if (error != 0) {
390             ReleaseMutex(floppy_lock);
391             return error;
392         }
393
394         memcpy(flp->cylinder_zero, flp->cylinder_buffer, 0x4800);
395
396         /*
397          * Must do this as we trashed the cached cylinder.
398          * Luckily we only ever need to do this once.
399          */
400         flp->stored_cylinder = cylinder;
401     }
402
403     PerformTransfer(flp->cylinder_zero + (512 * (sector - 1 + head * 18)), io, 512);
404 }
405
406 if (cylinder != flp->stored_cylinder) {
407     int error = FloppyDoCylinder(flp, cylinder);
408
409     if (error != 0) {
410         ReleaseMutex(floppy_lock);
411         return error;
412     }
413
414     PerformTransfer(flp->cylinder_buffer + (512 * (sector - 1 + head * 18)), io, 512);
415     flp->stored_cylinder = cylinder;
416 }
417
418 /*
419  * Read only read one sector, so we need to repeat the process if multiple
420  * sectors were requested. We could just do a larger copy above, but this is
421  * a bit simpler, as we don't need to worry about whether the entire request
422  * is on cylinder or not.
423  */
424 if (count > 0) {
425     ++lba;
426     --count;
427     goto next_sector;
428 }
429
430 ReleaseMutex(floppy_lock);
431 return 0;
432
433 static int ReadWrite(struct vnode* node, struct transfer* io) {
434     return FloppyIo(node->data, io);
435 }
436
437 static int Create(struct vnode* node, struct vnode** partition, const char* name, int, mode_t) {
438     AcquireMutex(floppy_lock, -1);
439     struct floppy_data* flp = node->data;
440     int res = DiskCreateHelper(&flp->partitions, partition, name);
441     ReleaseMutex(floppy_lock);
442     return res;
443 }
444
445 static int Follow(struct vnode* node, struct vnode** output, const char* name) {
446     AcquireMutex(floppy_lock, -1);
447     struct floppy_data* flp = node->data;
448     int res = DiskFollowHelper(&flp->partitions, output, name);
449     ReleaseMutex(floppy_lock);
450     return res;
451 }
452
453 static const struct vnode_operations dev_ops = {
454     .read = ReadWrite,
455     .write = ReadWrite,
456     .create = Create,
457     .follow = Follow,
458 };
459
460 void InitFloppy(void) {
461     floppy_lock = CreateMutex("floppy");
462
463     CreateThread(FloppyMotorControlThread, NULL, GetVas(), "flpmotor");
464
465     for (int i = 0; i < 1; ++i) {
466         struct vnode* node = CreateVnode(dev_ops, (struct stat) {
467             .st_mode = S_IFBLK | S_IRWXU | S_IRWXG | S_IRWXO,
468             .st_nlink = 1,
469             .st_blksize = 512,
470             .st_blocks = 2880,
471             .st_size = 512 * 2880
472         });
473
474         struct floppy_data* flp = AllocHeap(sizeof(struct floppy_data));
475         *flp = (struct floppy_data) {
476             .disk_num = i, .base = 0x3F0,
477             .stored_cylinder = -1, .got_cylinder_zero = false,
478             .cylinder_buffer = (uint8_t*) MapVirt(0, 0, CYLINDER_SIZE, VM_READ | VM_WRITE | VM_LOCK, NULL, 0);
479         };
480     }

```

```

481     .cylinder_zero = (uint8_t*) MapVirt(0, 0, CYLINDER_SIZE, VM_READ | VM_WRITE | VM_LOCK, NULL, 0);
482 };
483 node->data = flp;
484
485 RegisterIrqHandler(PIC_IRQ_BASE + 6, FloppyIrqHandler);
486 FloppyReset(flp);
487
488 InitDiskPartitionHelper(&flp->partitions);
489
490 AddVfsMount(node, GenerateNewRawDiskName(DISKUTIL_TYPE_FLOPPY));
491 CreateDiskPartitions(CreateOpenFile(node, 0, 0, true, true));
492 }
493

```

File: /arch/x86/dev/ide.c

```

1  #include <common.h>
2  #include <semaphore.h>
3  #include <log.h>
4  #include <thread.h>
5  #include <vfs.h>
6  #include <transfer.h>
7  #include <assert.h>
8  #include <errno.h>
9  #include <machine/portio.h>
10 #include <heap.h>
11 #include <virtual.h>
12 #include <stdlib.h>
13 #include <sys/stat.h>
14 #include <dirent.h>
15 #include <diskcache.h>
16 #include <irq.h>
17 #include <diskutil.h>
18
19 #define MAX_TRANSFER_SIZE (1024 * 16)
20
21 struct semaphore* ide_lock = NULL;
22
23 struct ide_data {
24     int disk_num;
25     unsigned int sector_size;
26     uint64_t total_num_sectors;
27     uint16_t* transfer_buffer;
28     size_t primary_base;
29     size_t primary_alternative;
30     size_t secondary_base;
31     size_t secondary_alternative;
32     size_t busmaster_base;
33
34     struct disk_partition_helper partitions;
35 };
36
37 int IdeCheckError(struct ide_data* ide) {
38     uint16_t base = ide->disk_num >= 2 ? ide->secondary_base : ide->primary_base;
39
40     uint8_t status = inb(base + 0x7);
41     if (status & 0x01) {
42         return EIO;
43     } else if (status & 0x20) {
44         return EIO;
45     } else if (!(status & 0x08)) {
46         return EIO;
47     }
48     return 0;
49 }
50
51 int IdePoll(struct ide_data* ide) {
52     uint16_t base = ide->disk_num >= 2 ? ide->secondary_base : ide->primary_base;
53     uint16_t alt_status_reg = ide->disk_num >= 2 ? ide->secondary_alternative : ide->primary_alternative;
54
55     /*
56      * Delay for a moment by reading the alternate status register.
57      */
58     for (int i = 0; i < 4; ++i) {
59         inb(alt_status_reg);
60     }
61
62     /*
63      * Wait for the device to not be busy. We have a timeout in case the
64      * device is faulty, we don't want to be in an endless loop and freeze
65      * the kernel.
66      */
67     int timeout = 0;
68     while (inb(base + 0x7) & 0x80) {
69         if (timeout > 975) {
70             SleepMilli(10);
71         }
72         if (timeout++ > 1000) {
73             return EIO;
74         }
75     }
76     return 0;
77 }
78
79 /*
80 * Read or write the primary ATA drive on the first controller. We use LBA28,
81 * so we are limited to a 28 bit sector number (i.e. disks up to 128GB in size)
82 */
83 static int IdeIo(struct ide_data* ide, struct transfer* io) {
84     EXACT_IRQL(IRQL_STANDARD);
85     int disk_num = ide->disk_num;
86
87     /*
88      * IDE devices do not contain an (accessible) disk buffer in PIO mode, as
89      * they transfer data through the IO ports. Hence we must read/write into
90      * this buffer first, and then move it safely to the destination.
91      * (we could use DMA instead, but PIO is simpler)
92      */
93     /* Allow up to 4KB sector sizes. Make sure there is enough room on the
94      * stack to handle this.
95      */
96     uint16_t* buffer = ide->transfer_buffer;
97
98     int sector = io->offset / ide->sector_size;
99     int count = io->length_remaining / ide->sector_size;
100
101     if (io->offset % ide->sector_size != 0) {
102         return EINVAL;
103     }
104     if (io->length_remaining % ide->sector_size != 0) {
105         return EINVAL;
106     }
107 }
108

```

```

111 if (count <= 0 || sector < 0 || sector > 0xFFFFFFFF || (uint64_t) sector + count >= ide->total_num_sectors) {
112     return EINVAL;
113 }
114
115 uint16_t base = disk_num >= 2 ? ide->secondary_base : ide->primary_base;
116 uint16_t dev_ctrl_reg = disk_num >= 2 ? ide->secondary_alternative : ide->primary_alternative;
117
118 int max_sectors_at_once = MAX_TRANSFER_SIZE / ide->sector_size;
119 if (max_sectors_at_once > 255) {
120     // hardware limitation
121     max_sectors_at_once = 255;
122 }
123
124 AcquireSemaphore(ide_lock, -1);
125
126 while (count > 0) {
127     int sectors_in_this_transfer = count > max_sectors_at_once ? max_sectors_at_once : count;
128
129     if (io->direction == TRANSFER_WRITE) {
130         PerformTransfer(buffer, io, ide->sector_size);
131     }
132
133     /*
134     * Send a whole heap of flags and the high 4 bits of the LBA to the controller.
135     */
136     outb(base + 0x6, 0xE0 | ((disk_num & 1) << 4) | ((sector >> 24) & 0xF));
137
138     /*
139     * Disable interrupts, we are going to use polling.
140     */
141     outb(dev_ctrl_reg, 2);
142
143     /*
144     * May not be needed, but it doesn't hurt to do it.
145     */
146     outb(base + 0x1, 0x00);
147
148     /*
149     * Send the number of sectors, and the sector's LBA.
150     */
151     outb(base + 0x2, sectors_in_this_transfer);
152     outb(base + 0x3, (sector >> 0) & 0xFF);
153     outb(base + 0x4, (sector >> 8) & 0xFF);
154     outb(base + 0x5, (sector >> 16) & 0xFF);
155
156     /*
157     * Send either the read or write command.
158     */
159     outb(base + 0x7, io->direction == TRANSFER_WRITE ? 0x30 : 0x20);
160
161     /*
162     * Wait for the data to be ready.
163     */
164     IdePoll(ide);
165
166     /*
167     * Read/write the data from/to the disk using ports.
168     */
169     if (io->direction == TRANSFER_WRITE) {
170         for (int c = 0; c < sectors_in_this_transfer; ++c) {
171             if (c != 0) {
172                 PerformTransfer(buffer, io, ide->sector_size);
173                 IdePoll(ide);
174             }
175
176             for (uint64_t i = 0; i < ide->sector_size / 2; ++i) {
177                 outw(base + 0x00, buffer[i]);
178             }
179
180             IdePoll(ide);
181
182             /*
183             * We need to flush the (hardware) disk cache if we are writing.
184             */
185             outb(base + 0x7, 0xE7);
186             IdePoll(ide);
187         }
188     } else {
189         int err = IdeCheckError(ide);
190         if (err) {
191             ReleaseSemaphore(ide_lock);
192             return err;
193         }
194
195         for (int c = 0; c < sectors_in_this_transfer; ++c) {
196             if (c != 0) {
197                 IdePoll(ide);
198             }
199
200             for (uint64_t i = 0; i < ide->sector_size / 2; ++i) {
201                 buffer[i] = inw(base + 0x00);
202             }
203
204             PerformTransfer(buffer, io, ide->sector_size);
205         }
206
207         /*
208         * Get ready for the next part of the transfer.
209         */
210         count -= sectors_in_this_transfer;
211         sector += sectors_in_this_transfer;
212     }
213
214     ReleaseSemaphore(ide_lock);
215     return 0;
216 }
217
218 static int IdeGetNumSectors(struct ide_data* ide) {
219     AcquireSemaphore(ide_lock, -1);
220
221     uint16_t base = ide->disk_num >= 2 ? ide->secondary_base : ide->primary_base;
222
223     /*
224     * Select the correct drive.
225     */
226     outb(base + 0x6, 0xE0 | ((ide->disk_num & 1) << 4));
227
228     /*
229     * Send the READ NATIVE MAX ADDRESS command, which will return the size
230     * of disk in sectors.
231     */
232     outb(base + 0x7, 0xF8);
233
234     IdePoll(ide);
235
236     /*
237     * The outputs are in the same registers we use to put the LBA
238     * when we read/write from the disk.
239     */
240     int sectors = 0;

```

```

241     sectors |= ((int) inb(base + 0x3));
242     sectors |= (((int) inb(base + 0x4)) << 8);
243     sectors |= (((int) inb(base + 0x5)) << 16);
244     sectors |= (((int) inb(base + 0x6) & 0xF) << 24);
245
246     ReleaseSemaphore(&ide_lock);
247
248     return sectors;
249 }
250
251 static int ReadWrite(struct vnode* node, struct transfer* io) {
252     return IdeIo(node->data, io);
253 }
254
255 static int Create(struct vnode* node, struct vnode** partition, const char* name, int, mode_t) {
256     AcquireSemaphore(&ide_lock, -1);
257     struct ide_data* ide = node->data;
258     int res = DiskCreateHelper(&ide->partitions, partition, name);
259     ReleaseSemaphore(&ide_lock);
260     return res;
261 }
262
263 static int Follow(struct vnode* node, struct vnode** output, const char* name) {
264     AcquireSemaphore(&ide_lock, -1);
265     struct ide_data* ide = node->data;
266     int res = DiskFollowHelper(&ide->partitions, output, name);
267     ReleaseSemaphore(&ide_lock);
268     return res;
269 }
270
271 static const struct vnode_operations dev_ops = {
272     .read = ReadWrite,
273     .write = ReadWrite,
274     .create = Create,
275     .follow = Follow,
276 };
277
278 void InitIde(void) {
279     ide_lock = CreateMutex("ide");
280
281     for (int i = 0; i < 1; ++i) {
282         struct ide_data* ide = AllocHeap(sizeof(struct ide_data));
283         *ide = (struct ide_data) {
284             .disk_num = i, .sector_size = 512, .busmaster_base = 0x0000,
285             .primary_base = 0x1F0, .primary_alternative = 0x3F6,
286             .secondary_base = 0x170, .secondary_alternative = 0x376,
287             .transfer_buffer = (uint16_t*) MapVirt(0, 0, MAX_TRANSFER_SIZE, VM_READ | VM_WRITE | VM_LOCK, NULL, 0);
288         };
289
290         ide->total_num_sectors = IdeGetNumSectors(ide);
291
292         struct vnode* node = CreateVnode(&dev_ops, (struct stat) {
293             .st_mode = S_IFBLK | S_IRWXU | S_IRWXG | S_IRWXO,
294             .st_nlink = 1,
295             .st_blksize = ide->sector_size,
296             .st_blocks = ide->total_num_sectors,
297             .st_size = ide->total_num_sectors * ide->sector_size,
298         });
299
300         InitDiskPartitionHelper(&ide->partitions);
301
302         node->data = ide;
303         AddVfsMount(node, GenerateNewRawDiskName(DISKUTIL_TYPE_FIXED));
304         CreateDiskPartitions(CreateDiskCache(CreateOpenFile(node, 0, 0, true, true)));
305     }
306 }

```

File: ./arch/x86/dev/pit.c

```

1
2 #include <machine/pic.h>
3 #include <machine/pit.h>
4 #include <machine/portio.h>
5 #include <arch.h>
6 #include <assert.h>
7 #include <common.h>
8 #include <timer.h>
9 #include <irq.h>
10 #include <log.h>
11
12 static uint64_t pit_hertz = 0;
13 static uint64_t pit_nanos = 0;
14
15 static int HandlePit(struct x86_regs*) {
16     ReceivedTimer(pit_nanos);
17     return 0;
18 }
19
20 void InitPit(int hertz) {
21     pit_hertz = (uint64_t) hertz;
22     pit_nanos = 1000000000ULL / pit_hertz;
23
24     int divisor = 1193180 / hertz;
25     outb(0x43, 0x36);
26     outb(0x40, divisor & 0xFF);
27     outb(0x40, divisor >> 8);
28
29     RegisterIrqHandler(PIC_IRQ_BASE + 0, HandlePit);
30 }

```

File: ./arch/x86/lowlevel/trap.s

```

1
2 global isr0
3 global isr1
4 global isr2
5 global isr3
6 global isr4
7 global isr5
8 global isr6
9 global isr7
10 global isr8
11 global isr9
12 global isr10
13 global isr11
14 global isr12
15 global isr13
16 global isr14
17 global isr15
18 global isr16
19 global isr17
20 global isr18
21 global isr19

```

```

22 global isr20
23 global isr21
24 global isr96
25 global irq0
26 global irq1
27 global irq2
28 global irq3
29 global irq4
30 global irq5
31 global irq6
32 global irq7
33 global irq8
34 global irq9
35 global irq10
36 global irq11
37 global irq12
38 global irq13
39 global irq14
40 global irq15
41
42
43 ; We don't need to disable interrupts - they are automatically disabled when
44 ; the interrupt comes in
45
46 isr0:
47     push 0
48     push 0
49     jmp int_common_handler
50
51 isr1:
52     push byte 0
53     push byte 1
54     jmp int_common_handler
55
56 isr2:
57     push byte 0
58     push byte 2
59     jmp int_common_handler
60
61 isr3:
62     push byte 0
63     push byte 3
64     jmp int_common_handler
65
66 isr4:
67     push byte 0
68     push byte 4
69     jmp int_common_handler
70
71 isr5:
72     push byte 0
73     push byte 5
74     jmp int_common_handler
75
76 isr6:
77     push byte 0
78     push byte 6
79     jmp int_common_handler
80
81 isr7:
82     push byte 0
83     push byte 7
84     jmp int_common_handler
85
86 isr8:
87     push byte 8
88     jmp int_common_handler
89
90 isr9:
91     push byte 0
92     push byte 9
93     jmp int_common_handler
94
95 isr10:
96     push byte 10
97     jmp int_common_handler
98
99 isr11:
100    push byte 11
101    jmp int_common_handler
102
103 isr12:
104    push byte 12
105    jmp int_common_handler
106
107 isr13:
108    push byte 13
109    jmp int_common_handler
110
111 isr14:
112    push byte 14
113    jmp int_common_handler
114
115 isr15:
116    push byte 0
117    push byte 15
118    jmp int_common_handler
119
120 isr16:
121    push byte 0
122    push byte 16
123    jmp int_common_handler
124
125 isr17:
126    push byte 17
127    jmp int_common_handler
128
129 isr18:
130    push byte 0
131    push byte 18
132    jmp int_common_handler
133
134 isr19:
135    push byte 0
136    push byte 19
137    jmp int_common_handler
138
139 isr20:
140    push byte 0
141    push byte 20
142    jmp int_common_handler
143
144 isr21:
145    push byte 0
146    push byte 21
147    jmp int_common_handler
148
149
150 ; This is our system call handler
151 isr96:

```

```

152     push byte 0
153     push 96
154     jmp int_common_handler
155
156
157
158 ; Note that in the PIC setup, we remap our IRQs so they start at 32
159 ; That is, IRQ0 is actually ISR32, etc. up to IRQ15 which is ISR47
160 ; This is so they don't clash with the exceptions above, which are
161 ; not re-mappable.
162 irq0:
163     push byte 0
164     push byte 32
165     jmp int_common_handler
166
167 irq1:
168     push byte 0
169     push byte 33
170     jmp int_common_handler
171
172 irq2:
173     push byte 0
174     push byte 34
175     jmp int_common_handler
176
177 irq3:
178     push byte 0
179     push byte 35
180     jmp int_common_handler
181
182 irq4:
183     push byte 0
184     push byte 36
185     jmp int_common_handler
186
187 irq5:
188     push byte 0
189     push byte 37
190     jmp int_common_handler
191
192 irq6:
193     push byte 0
194     push byte 38
195     jmp int_common_handler
196
197 irq7:
198     push byte 0
199     push byte 39
200     jmp int_common_handler
201
202 irq8:
203     push byte 0
204     push byte 40
205     jmp int_common_handler
206
207 irq9:
208     push byte 0
209     push byte 41
210     jmp int_common_handler
211
212 irq10:
213     push byte 0
214     push byte 42
215     jmp int_common_handler
216
217 irq11:
218     push byte 0
219     push byte 43
220     jmp int_common_handler
221
222 irq12:
223     push byte 0
224     push byte 44
225     jmp int_common_handler
226
227 irq13:
228     push byte 0
229     push byte 45
230     jmp int_common_handler
231
232 irq14:
233     push byte 0
234     push byte 46
235     jmp int_common_handler
236
237 irq15:
238     push byte 0
239     push byte 47
240     jmp int_common_handler
241
242
243 ; Our common interrupt handler
244 extern x86HandleInterrupt
245 int_common_handler:
246     ; Save the registers and segments
247     pushad
248     push ds
249     push es
250     push fs
251     push gs
252
253     ; Ensure we have kernel segments and not user segments
254     mov ax, 0x10
255     mov ds, ax
256     mov es, ax
257     ; We are going to use FS to store the CPU number, so don't overwrite it
258     ; Kernel has no use for GS, so don't bother setting it
259
260     ; Allow nested interrupts
261     sti
262
263     ; Push a pointer to the registers to the kernel handler
264     push esp
265     cld
266     call x86HandleInterrupt
267
268     ; Restore registers
269     add esp, 4
270     pop gs
271     pop fs
272     pop es
273     pop ds
274     popad
275
276     ; Skip over the error code and interrupt number (we can't pop them
277     ; anywhere, as the registers have already been restored)
278     add esp, 8
279
280     ; Return from the interrupt - also restores the stack and the flags
281     iretd

```

File: `/arch/x86/lowlevel/ldt.s`

```
1
2 global x86LoadIdt
3 x86LoadIdt:
4 ; The address of the IDTR is passed in as an argument
5 mov eax, [esp + 4]
6 ldt [eax]
7
8 ret
```

File: `/arch/x86/lowlevel/gdt.c`

```
1
2 #include <common.h>
3 #include <cpu.h>
4 #include <machine/gdt.h>
5
6 extern void x86LoadGdt(size_t addr);
7
8 static struct gdt_entry x86CreateGdtEntry(size_t base, size_t limit, uint8_t access, uint8_t granularity)
9 {
10     struct gdt_entry entry;
11
12     entry.base_low = base & 0xFFFF;
13     entry.base_middle = (base >> 16) & 0xFF;
14     entry.base_high = (base >> 24) & 0xFF;
15     entry.limit_low = limit & 0xFFFF;
16     entry.flags_and_limit_high = (limit >> 16) & 0xF;
17     entry.flags_and_limit_high |= (granularity & 0xF) << 4;
18     entry.access = access;
19
20     return entry;
21 }
22
23 void x86InitGdt(void)
24 {
25     platform_cpu_data_t* cpu_data = GetCpu()->platform_specific;
26
27     cpu_data->gdt[0] = x86CreateGdtEntry(0, 0, 0, 0); // null segment
28     cpu_data->gdt[1] = x86CreateGdtEntry(0, 0xFFFFFFFF, 0x9A, 0xC); // kernel code
29     cpu_data->gdt[2] = x86CreateGdtEntry(0, 0xFFFFFFFF, 0x92, 0xC); // kernel data
30     cpu_data->gdt[3] = x86CreateGdtEntry(0, 0xFFFFFFFF, 0xFA, 0xC); // user code
31     cpu_data->gdt[4] = x86CreateGdtEntry(0, 0xFFFFFFFF, 0xF2, 0xC); // user data
32
33     cpu_data->gdt.size = sizeof(cpu_data->gdt) - 1;
34     cpu_data->gdt.location = (size_t) &cpu_data->gdt;
35
36     x86LoadGdt((size_t) &cpu_data->gdt);
37 }
38
39 /*
40 * Adds a Task State Segment (TSS) entry into the GDT. This allows the TSS to be
41 * used to switch from user mode to kernel mode.
42 *
43 * Returns the selector used in the GDT.
44 */
45
46 uint16_t x86AddTssToGdt(struct tss* tss)
47 {
48     platform_cpu_data_t* cpu_data = GetCpu()->platform_specific;
49     cpu_data->gdt[5] = x86CreateGdtEntry((size_t) tss, sizeof(struct tss), 0x89, 0x0);
50     return 5 * 0x8;
51 }
```

File: `/arch/x86/lowlevel/tss.c`

```
1
2 #include <heap.h>
3 #include <machine/tss.h>
4 #include <arch.h>
5 #include <machine/gdt.h>
6 #include <cpu.h>
7
8 extern void x86LoadTss(size_t selector);
9
10 void x86InitTss(void) {
11     platform_cpu_data_t* cpu_data = GetCpu()->platform_specific;
12     cpu_data->tss = (struct tss*) AllocHeap(sizeof(struct tss));
13
14     cpu_data->tss->link = 0x10;
15     cpu_data->tss->esp0 = 0;
16     cpu_data->tss->ss0 = 0x10;
17     cpu_data->tss->iopb = sizeof(struct tss);
18
19     uint16_t selector = x86AddTssToGdt(cpu_data->tss);
20     x86LoadTss(selector);
21 }
```

File: `/arch/x86/lowlevel/misc.s`

```

1
2 global ArchReadTimestamp
3 global ArchEnableInterrupts
4 global ArchDisableInterrupts
5 global ArchStallProcessor
6 global ArchFlushTlb
7 global x86GetCr2
8 global x86AreCpusOn
9 global ArchGetCurrentCpuIndex
10
11 ArchGetCurrentCpuIndex:
12     ; Needs to be something that can't be modified by user code (e.g. a debug register).
13     mov eax, dr3
14     ret
15
16 ArchReadTimestamp:
17     rdtsc
18     ret
19
20 ArchEnableInterrupts:
21     sti
22     ret
23
24 ArchDisableInterrupts:
25     cli
26     ret
27
28 ArchStallProcessor:
29     hlt
30     ret
31
32 x86AreCpusOn:
33     pushf
34     pop eax
35     and eax, 0x200
36     shr eax, 9
37     ret

```

File: `/arch/x86/lowlevel/idt.c`


```

1
2 #include <common.h>
3 #include <cpu.h>
4
5 /*
6  * x86/lowlevel/idt.c - Interrupt Descriptor Table
7  *
8  * The interrupt descriptor table (IDT) is essentially a lookup table for where
9  * the CPU should jump to when an interrupt is received.
10 */
11
12 extern void x86LoadIdt(size_t addr);
13
14 /*
15  * Our trap handlers, defined in lowlevel/trap.s, which will be called
16  * when an interrupt occurs.
17 */
18 extern void isr0();
19 extern void isr1();
20 extern void isr2();
21 extern void isr3();
22 extern void isr4();
23 extern void isr5();
24 extern void isr6();
25 extern void isr7();
26 extern void isr8();
27 extern void isr9();
28 extern void isr10();
29 extern void isr11();
30 extern void isr12();
31 extern void isr13();
32 extern void isr14();
33 extern void isr15();
34 extern void isr16();
35 extern void isr17();
36 extern void isr18();
37 extern void isr19();
38 extern void isr20();
39 extern void isr21();
40 extern void isr96();
41 extern void irq0();
42 extern void irq1();
43 extern void irq2();
44 extern void irq3();
45 extern void irq4();
46 extern void irq5();
47 extern void irq6();
48 extern void irq7();
49 extern void irq8();
50 extern void irq9();
51 extern void irq10();
52 extern void irq11();
53 extern void irq12();
54 extern void irq13();
55 extern void irq14();
56 extern void irq15();
57
58 /*
59  * Fill in an entry in the IDT. There are a number of 'types' of interrupt, determining
60  * whether interrupts are disabled automatically before calling the handler, whether
61  * it is a 32-bit or 16-bit entry, and whether user mode can invoke the interrupt manually.
62 */
63 static void x86SetIdtEntry(int num, size_t isr_addr, uint8_t type)
64 {
65     platform_cpu_data_t* cpu_data = GetCpu()->platform_specific;
66
67     cpu_data->idt[num].isr_offset_low = (isr_addr & 0xFFFF);
68     cpu_data->idt[num].isr_offset_high = (isr_addr >> 16) & 0xFFFF;
69     cpu_data->idt[num].segment_selector = 0x08;
70     cpu_data->idt[num].reserved = 0;
71     cpu_data->idt[num].type = type;
72 }
73
74 /*
75  * Initialise the IDT. After this has occurred, interrupts may be enabled.
76 */
77 void x86InitIdt(void)
78 {
79     platform_cpu_data_t* cpu_data = GetCpu()->platform_specific;
80
81     void (*const isrs[])() = {
82         isr0, isr1, isr2, isr3, isr4, isr5, isr6, isr7,
83         isr8, isr9, isr10, isr11, isr12, isr13, isr14, isr15,
84         isr16, isr17, isr18, isr19, isr20, isr21,
85     };
86
87     void (*const irq_handlers[])() = {
88         irq0, irq1, irq2, irq3, irq4, irq5, irq6, irq7,
89         irq8, irq9, irq10, irq11, irq12, irq13, irq14, irq15
90     };
91
92     /*
93      * Install handlers for CPU exceptions (e.g. for page faults, divide-by-zero, etc.).
94      */
95     for (int i = 0; i < 21; ++i) {
96         x86SetIdtEntry(i, (size_t) isrs[i], 0x8E);
97     }
98
99     /*
100      * Install handlers for IRQs (hardware interrupts, e.g. keyboard, system timer).
101      */
102     for (int i = 0; i < 16; ++i) {
103         x86SetIdtEntry(i + 32, (size_t) irq_handlers[i], 0x8E);
104     }
105
106     /*
107      * Install our system call handler. Note that the flag byte is 0xEE instead of 0x8E,
108      * this allows user code to directly invoke this interrupt.
109      */
110     x86SetIdtEntry(96, (size_t) isr96, 0xEE);
111
112     cpu_data->idtr.location = (size_t) cpu_data->idt;
113     cpu_data->idtr.size = sizeof(cpu_data->idt) - 1;
114
115     x86LoadIdt((size_t) cpu_data->idtr);
116 }

```

File: ./arch/x86/lowlevel/tss.s

```

1
2 ;
3 ;
4 ; x86/lowlevel/tss.s - Task State Segment
5 ;
6 ; Like with the GDT and IDT, we need assembly to load the TSS using the
7 ; special instruction 'ltr'.
8 ;
9
10 extern x86LoadTss
11
12 x86LoadTss:
13     mov eax, [esp + 4]
14     ltr ax
15     ret

```

File: /arch/x86/lowlevel/gdt.s

```

1
2
3 global x86LoadGdt
4 x86LoadGdt:
5 ; The address of the GDTR is passed in as an argument
6 mov eax, [esp + 4]
7 lgdt [eax]
8
9 ; We now need to reload CS using a far jump...
10 jmp 0x08:.reloadSegments
11
12 .reloadSegments:
13 ; And all of the other segments by loading them
14 mov ax, 0x10
15 mov ds, ax
16 mov es, ax
17 ; Kernel doesn't use gs/fs
18 mov ss, ax
19
20 ret

```

File: /arch/x86/application.ld

```

ENTRY(_start)
OUTPUT_FORMAT("elf32-i386")

SECTIONS
{
    . = 0x10000000;

    .text ALIGN(4096) : AT (ADDR (.text))
    {
        *(.text)
        *(.rodata)
        *(.symtab)
        *(.strtab)
    }

    .data ALIGN(4096) : AT (ADDR (.data))
    {
        *(.data)
    }

    .bss ALIGN(4096) : AT (ADDR (.bss))
    {
        *(COMMON)
        *(.bss)
        *(.bootstrap_stack)
    }

    /DISCARD/ :
    {
        *(.comment)
    }
}

```

File: /arch/x86/thread/spinlock.s

```

1
2 ;
3 ;
4 ; x86/thread/spinlock.s - Spinlocks
5 ;
6 ; Implement spinlocks in assembly so we can guarantee that they are
7 ; atomic.
8 ;
9 ;
10 ;
11
12 global ArchSpinlockAcquire
13 global ArchSpinlockRelease
14
15 ArchSpinlockAcquire:
16 ; The address of the lock is passed in as an argument
17 mov eax, [esp + 4]
18
19 .try_acquire:
20 ; Try to acquire the lock
21 lock bts dword [eax], 0
22 jc .spin_wait
23
24 ret
25
26 .spin_wait:
27 ; Lock was not acquired, so do the 'spin' part of spinlock
28
29 ; Hint to the CPU that we are spinning
30 pause
31
32 ; No point trying to acquire it until it is free
33 test dword [eax], 1
34 jnz .spin_wait
35
36 ; Now that it is free, we can attempt to atomically acquire it again
37 jmp .try_acquire
38
39
40 ArchSpinlockRelease:
41 ; The address of the lock is passed in as an argument
42 mov eax, [esp + 4]
43 lock btr dword [eax], 0
44 ret

```

File: /arch/x86/thread/usermode.s

```

1
2 global ArchSwitchToUsermode
3
4 ; This is only called the first time we want to switch a given
5 ; thread into usermode. In all other cases the switch will occur
6 ; back through an interrupt handler (e.g. after a system call completes)
7
8 ArchSwitchToUsermode:
9 ; Takes in an address to a usermode address to start execution
10 mov ebx, [esp + 4]
11
12 ; And a user stack pointer
13 mov ecx, [esp + 8]
14
15 ; And initial argument. For x86, we'll just store this in EDI, and the
16 ; program can just read that value.
17 mov edi, [esp + 12]
18
19 ; Usermode data segment
20 mov ax, 0x23
21 mov ds, ax
22 mov es, ax
23 mov fs, ax
24 mov gs, ax
25
26 ; We need to push the current stack as the usermode part will use it too
27 push 0x23 ; Usermode stack segment
28 push ecx ; Usermode stack pointer
29 push 0x202 ; Flags
30 push 0x1B ; Usermode code segment
31 push ebx ; Usermode entry point
32
33 ; Pop all of that off the stack and go to usermode.
34 iretd

```

File: /arch/x86/thread/switch.s

```

1 global ArchPrepareStack
2 global ArchSwitchThread
3
4 extern ThreadInitialisationHandler
5 extern GetCpu
6
7 ArchPrepareStack:
8 ; We need to put 5 things on the stack - dummy values for EBX, ESI,
9 ; EDI and EBP, as well as the address of thread_startup_handler
10 ;
11 ; This is because these get popped off in arch_switch_thread
12
13 ; Grab the address of the new thread's stack from our stack (it was
14 ; passed in as an argument)
15 mov eax, [esp + 4]
16
17 ; We need to get to the bottom position, and we also need to return that
18 ; address in EAX so it can be put into the struct, so it makes sense to modify it.
19 sub eax, 20
20
21 ; This is where the address of arch_switch_thread needs to go.
22 ; +0 is where EBP is, +4 is EDI, +8 for ESI, +12 for EBX,
23 ; and so +16 for the return value,
24 ; (see the start of arch_switch_thread for where these get pushed)
25 mov [eax + 16], dword ThreadInitialisationHandler
26
27 ret
28
29 ArchSwitchThread:
30 ; The old and new threads are passed in on the stack as arguments, in that order.
31
32 ; The calling convention we use already saves EAX, ECX and EDX whenever a
33 ; function (e.g. ArchSwitchThread) is called. Therefore, we only need to save the other four.
34 push ebx
35 push esi
36 push edi
37 push ebp
38
39 ; We are now free to trash the general purpose registers (except ESP),
40 ; so we can now load the current task using the argument.
41
42 ; First we have to save the old stack pointer. The old thread was the first
43 ; argument, and we just pushed 4 things to the stack. The first argument gets
44 ; pushed last, so read back 5 places. Also load the new thread's address in.
45
46 mov edi, [esp + (4 + 1) * 4] ; edi = old_thread
47 mov esi, [esp + (4 + 2) * 4] ; esi = new_thread
48
49 ; The second entry in a thread structure is guaranteed to be the stack pointer.
50 ; Save our stack there.
51 mov [edi + 4], esp ; old_thread->stack_pointer = esp
52
53 ; Now we can load the new thread's stack pointer.
54 mov esp, [esi + 4] ; esp = new_thread->stack_pointer
55
56 ; ESI is callee-saved, so no need to do anything here. We only need ESI and ESP
57 ; at this point, so it's all good.
58
59 call GetCpu ; eax = GetCpu()
60
61 ; The top of the kernel stack (which needs to go in the TSS for
62 ; user to kernel switches), is the first entry in new_thread.
63 mov ebx, [esi] ; ebx = new_thread->kernel_stack_top
64
65 ; The third entry in current_cpu is a pointer to CPU specific data.
66 mov ecx, [eax + 8] ; ecx = GetCpu()->platform_specific
67
68 ; The first entry in the CPU specific data is the TSS pointer
69 mov edx, [ecx + 0] ; edx = GetCpu()->platform_specific->tss
70
71 ; Load the TSS's ESP0 with the new thread's stack
72 mov [edx + 4], ebx
73
74 ; Now we have the new thread's stack, we can just pop off the state
75 ; that would have been pushed when it was switched out.
76 pop ebp
77 pop edi
78 pop esi
79 pop ebx
80
81 ret

```

File: /arch/x86/progload.ld

```

ENTRY(_start)
OUTPUT_FORMAT("binary")

SECTIONS
{
    . = 0xBFC00000;

    .text ALIGN(4096) : AT (ADDR (.text))
    {
        *(.text)
        *(.rodata)
        *(.symtab)
        *(.strtab)
    }

    .data ALIGN(4096) : AT (ADDR (.data))
    {
        *(.data)
    }

    .bss ALIGN(4096) : AT (ADDR (.bss))
    {
        *(COMMON)
        *(.bss)
        *(.bootstrap_stack)
    }

    .fake : { . = . + SIZEOF(.bss); }

/DISCARD/ :
{
    *(.comment)
}

```

File: /arch/x86/linker.ld

```
ENTRY(_start)
OUTPUT_FORMAT("elf32-i386")

SECTIONS
{
    . = 1M;

    _kernel_start = .;
    .multiboot.data : {
        *(.multiboot.data)
    }

    .multiboot.text : {
        *(.multiboot.text)
    }

    . += 0xC0000000;

    .text ALIGN(4096) : AT (ADDR (.text) - 0xC0000000)
    {
        *(.text.text.*)
        *(.ctors)
        *(.dtors)
    }

    .rodata ALIGN(4096) : AT (ADDR (.rodata) - 0xC0000000)
    {
        *(.rodata)
    }

    . = ALIGN(4096);
    _start_pageablek_section = .;

    .pageablek ALIGN(4096) : AT (ADDR (.data) + SIZEOF(.data) - 0xC0000000)
    {
        *(.pageablektext)
        *(.pageablekdata)
    }
    _end_pageablek_section = .;

    .data ALIGN(4096) : AT (ADDR (.data) - 0xC0000000)
    {
        *(.data)
    }

    .bss ALIGN(4096) : AT (ADDR (.bss) - 0xC0000000)
    {
        *(COMMON)
        *(.bss)
        *(.bootstrap_stack)
    }

    _kernel_end = .;

    /DISCARD/ :
    {
        *(.comment)
    }
}
```

File: /arch/x86/elf/elf.c

```
1  #include <common.h>
2  #include <errno.h>
3  #include <log.h>
4  #include <vfs.h>
5  #include <virtual.h>
6  #include <voidptr.h>
7  #include <string.h>
8  #include <driver.h>
9  #include <sys/stat.h>
10 #include <assert.h>
11 #include <heap.h>
12 #include <irq.h>
13 #include <stdlib.h>
14 #include <machine/elf.h>
15 #include <panic.h>
16
17 static bool IsElfValid(struct Elf32_Ehdr* header) {
18     if (header->e_ident[EI_MAG0] != ELF_MAG0) return false;
19     if (header->e_ident[EI_MAG1] != ELF_MAG1) return false;
20     if (header->e_ident[EI_MAG2] != ELF_MAG2) return false;
21     if (header->e_ident[EI_MAG3] != ELF_MAG3) return false;
22
23     /* TODO: check for other things, such as the platform, etc. */
24
25     return true;
26 }
27
28 static size_t ElfGetSizeOfImageIncludingBss(void* data) {
29     struct Elf32_Ehdr* elf_header = (struct Elf32_Ehdr*) data;
30     struct Elf32_Phdr* prog_headers = (struct Elf32_Phdr*) AddVoidPtr(data, elf_header->e_phoff);
31
32     size_t base_point = 0xD00000000;
33     size_t total_size = 0;
34
35     for (int i = 0; i < elf_header->e_phnum; ++i) {
36         struct Elf32_Phdr* prog_header = prog_headers + i;
37
38         size_t address = prog_header->p_vaddr;
39         size_t size = prog_header->p_filesz;
40         size_t num_zero_bytes = prog_header->p_memsz - size;
41         size_t type = prog_header->p_type;
42
43         if (type == PHT_LOAD) {
44             if (address - base_point + size + num_zero_bytes >= total_size) {
45                 total_size = address - base_point + size + num_zero_bytes;
46             }
47         }
48     }
49
50     return (total_size + ARCH_PAGE_SIZE - 1) & ~(ARCH_PAGE_SIZE - 1);
}
```

```

51 }
52
53 static int ElfLoadProgramHeaders(void* data, size_t relocation_point, struct open_file* file) {
54     struct Elf32_Ehdr* elf_header = (struct Elf32_Ehdr*) data;
55     struct Elf32_Phdr* prog_headers = (struct Elf32_Phdr*) AddVoidPtr(data, elf_header->e_phoff);
56
57     size_t base_point = 0x00000000;
58
59     for (int i = 0; i < elf_header->e_phnum; ++i) {
60         struct Elf32_Phdr* prog_header = prog_headers + i;
61
62         size_t address = prog_header->p_vaddr;
63         size_t offset = prog_header->p_offset;
64         size_t size = prog_header->p_filesz;
65         size_t type = prog_header->p_type;
66         uint32_t flags = prog_header->p_flags;
67         size_t num_zero_bytes = prog_header->p_memsz - size;
68
69         if (type == PHT_LOAD) {
70             size_t addr = address + relocation_point - base_point;
71             size_t remainder = size % (ARCH_PAGE_SIZE - 1);
72
73             int page_flags = 0;
74             if (flags & PF_X) page_flags |= VM_EXEC;
75             if (flags & PF_W) page_flags |= VM_WRITE;
76             if (flags & PF_R) page_flags |= VM_READ;
77
78             /*
79              * We don't actually want to write to the executable file, so we must just copy to the page as normal
80              * instead of using a file-backed page.
81              */
82             if (flags & PF_W) {
83                 size_t pages = (size + num_zero_bytes + (ARCH_PAGE_SIZE - 1)) / ARCH_PAGE_SIZE;
84
85                 for (size_t i = 0; i < pages; ++i) {
86                     SetVirtPermissions(addr + i * ARCH_PAGE_SIZE, page_flags, (VM_READ | VM_WRITE | VM_EXEC) & ~page_flags);
87                 }
88
89                 memcpy((void*) addr, (const void*) AddVoidPtr(data, offset), size);
90
91             } else {
92                 size_t pages = (size - remainder) / ARCH_PAGE_SIZE;
93
94                 if (addr % (ARCH_PAGE_SIZE - 1)) {
95                     return EINVAL;
96                 }
97
98                 if (pages > 0) {
99                     LogWriteSerial("doing the little fiddly thing...\n");
100                     UnmapVirt(addr, pages * ARCH_PAGE_SIZE);
101                     size_t v = MapVirt(relocation_point, addr, pages * ARCH_PAGE_SIZE, VM_RELOCATABLE | VM_FILE | page_flags, file, offset);
102                     if (v != addr) {
103                         return ENOMEM;
104                     }
105                 }
106
107                 if (remainder > 0) {
108                     SetVirtPermissions(addr + pages * ARCH_PAGE_SIZE, page_flags | VM_WRITE, (VM_READ | VM_EXEC) & ~page_flags);
109                     memcpy((void*) AddVoidPtr(addr, pages * ARCH_PAGE_SIZE), (const void*) AddVoidPtr(data, offset + pages * ARCH_PAGE_SIZE), remainder);
110                     SetVirtPermissions(addr + pages * ARCH_PAGE_SIZE, 0, VM_WRITE);
111                 }
112             }
113         }
114     }
115     return 0;
116 }
117
118 static char* ElfLookupString(void* data, int offset) {
119     struct Elf32_Ehdr* elf_header = (struct Elf32_Ehdr*) data;
120
121     if (elf_header->e_shstrndx == SHN_UNDEF) {
122         return NULL;
123     }
124
125     struct Elf32_Shdr* sect_headers = (struct Elf32_Shdr*) AddVoidPtr(data, elf_header->e_shoff);
126
127     char* string_table = (char*) AddVoidPtr(data, sect_headers[elf_header->e_shstrndx].sh_offset);
128     if (string_table == NULL) {
129         return NULL;
130     }
131
132     return string_table + offset;
133 }
134
135 static size_t ElfGetSymbolValue(void* data, int table, size_t index, bool* error, size_t relocation_point, size_t base_address) {
136     *error = false;
137
138     if (table == SHN_UNDEF || index == SHN_UNDEF) {
139         *error = true;
140         return 0;
141     }
142
143     struct Elf32_Ehdr* elf_header = (struct Elf32_Ehdr*) data;
144     struct Elf32_Shdr* sect_headers = (struct Elf32_Shdr*) AddVoidPtr(data, elf_header->e_shoff);
145     struct Elf32_Shdr* symbol_table = sect_headers + table;
146     struct Elf32_Shdr* string_table = sect_headers + symbol_table->sh_link;
147
148     size_t num_symbol_table_entries = symbol_table->sh_size / symbol_table->sh_entsize;
149     if (index >= num_symbol_table_entries) {
150         *error = true;
151         return 0;
152     }
153
154     struct Elf32_Sym* symbol = ((struct Elf32_Sym*) AddVoidPtr(data, symbol_table->sh_offset)) + index;
155
156     if (symbol->st_shndx == SHN_UNDEF) {
157         const char* name = (const char*) AddVoidPtr(data, string_table->sh_offset + symbol->st_name);
158
159         size_t target = GetSymbolAddress(name);
160         if (target == 0) {
161             if (!(ELF32_ST_BIND(symbol->st_info) & STB_WEAK)) {
162                 *error = true;
163             }
164             return 0;
165         }
166     } else {
167         return target;
168     }
169
170     if (symbol->st_shndx == SHN_ABS) {
171         return symbol->st_value;
172     }
173     if (symbol->st_shndx == SHN_REL) {
174         return symbol->st_value + (relocation_point - base_address);
175     }
176 }
177
178 static bool ElfPerformRelocation(void* data, size_t relocation_point, struct Elf32_Shdr* section, struct Elf32_Rel* relocation_table, struct quick_relocation_tab
179
180     size_t base_address = 0x00000000;

```

```

181 size_t addr = (size_t) relocation_point - base_address + relocation_table->r_offset;
182 size_t* ref = (size_t*) addr;
183
184 int symbolValue = 0;
185 if (ELF32_R_SYM(relocation_table->r_info) != SHN_UNDEF) {
186     bool error = false;
187     symbolValue = ElfGetSymbolValue(data, section->sh_link, ELF32_R_SYM(relocation_table->r_info), &error, relocation_point, base_address);
188     if (error) {
189         return false;
190     }
191 }
192
193 bool needs_write_low = (GetVirtPermissions(addr) & VM_WRITE) == 0;
194 bool needs_write_high = (GetVirtPermissions(addr + sizeof(size_t) - 1) & VM_WRITE) == 0;
195
196 if (needs_write_low) {
197     SetVirtPermissions(addr, VM_WRITE, 0);
198 }
199 if (needs_write_high) {
200     SetVirtPermissions(addr + sizeof(size_t) - 1, VM_WRITE, 0);
201 }
202
203 int type = ELF32_R_TYPE(relocation_table->r_info);
204 bool success = true;
205 size_t val = 0;
206
207 if (type == R_386_32) {
208     val = DO_386_32(symbolValue, *ref);
209 }
210 else if (type == R_386_PC32) {
211     val = DO_386_PC32(symbolValue, *ref, addr);
212 }
213 else if (type == R_386_RELATIVE) {
214     val = DO_386_RELATIVE((relocation_point - base_address), *ref);
215 }
216 else {
217     LogWriteSerial("some whacko type...\n");
218     success = false;
219 }
220
221 *ref = val;
222 LogWriteSerial("relocating 0x%X -> 0x%X\n", addr, val);
223 AddToQuickRelocationTable(table, addr, val);
224
225 if (needs_write_low) {
226     SetVirtPermissions(addr, 0, VM_WRITE);
227 }
228 if (needs_write_high) {
229     SetVirtPermissions(addr + sizeof(size_t) - 1, 0, VM_WRITE);
230 }
231 return success;
232 }
233
234 static bool ElfPerformRelocations(void* data, size_t relocation_point, struct quick_relocation_table** table) {
235     struct Elf32_Ehdr* elf_header = (struct Elf32_Ehdr*) data;
236     struct Elf32_Shdr* sect_headers = (struct Elf32_Shdr*) AddVoidPtr(data, elf_header->e_shoff);
237
238     for (int i = 0; i < elf_header->e_shnum; ++i) {
239         struct Elf32_Shdr* section = sect_headers + i;
240
241         if (section->sh_type == SHT_REL) {
242             struct Elf32_Rel* relocation_tables = (struct Elf32_Rel*) AddVoidPtr(data, section->sh_offset);
243             int count = section->sh_size / section->sh_entsize;
244
245             if (strcmp(ElfLookupString(data, section->sh_name), ".rel.dyn")) {
246                 continue;
247             }
248
249             *table = CreateQuickRelocationTable(count);
250
251             for (int index = 0; index < count; ++index) {
252                 bool success = ElfPerformRelocation(data, relocation_point, section, relocation_tables + index, *table);
253                 if (!success) {
254                     LogWriteSerial("failed to do a relocation!! (%d)\n", index);
255                     return false;
256                 }
257             }
258
259             SortQuickRelocationTable(*table);
260
261         } else if (section->sh_type == SHT_RELA) {
262             LogDeveloperWarning("[ElfPerformRelocations]: unsupported section type: SHT_RELA\n");
263             return false;
264         }
265     }
266 }
267
268 return true;
269 }
270
271 static int ElfLoad(void* data, size_t* relocation_point, struct open_file* file, struct quick_relocation_table** table) {
272     MAX_IRQL(IRQL_PAGE_FAULT);
273
274     struct Elf32_Ehdr* elf_header = (struct Elf32_Ehdr*) data;
275
276     if (!IsValidElfHeader(elf_header)) {
277         return EINVAL;
278     }
279
280     /*
281     * To load a driver, we need the section headers.
282     */
283     if (elf_header->e_shnum == 0) {
284         return EINVAL;
285     }
286
287     /*
288     * We always need the program headers.
289     */
290     if (elf_header->e_phnum == 0) {
291         return EINVAL;
292     }
293
294     /*
295     * Load into memory.
296     */
297     size_t size = ElfGetSizeOfImageIncludingBss(data);
298
299     *relocation_point = MapVirt(0, 0, size, VM_READ, NULL, 0);
300     LogWriteSerial("RELOCATION POINT AT 0x%X\n", *relocation_point);
301     ElfLoadProgramHeaders(data, *relocation_point, file);
302
303     bool success = ElfPerformRelocations(data, *relocation_point, table);
304     if (success) {
305         return 0;
306     } else {
307         return EINVAL;
308     }
309 }
310

```

```

311
312 int ArchLoadDriver(size_t* relocation_point, struct open_file* file, struct quick_relocation_table** table) {
313     EXACT_IQRL(IQRL_STANDARD);
314
315     off_t file_size = file->node->stat.st_size;
316     size_t file_rgn = MapVirt(0, 0, file_size, VM_READ | VM_FILE, file, 0);
317     int res = ElfLoad((void*) file_rgn, relocation_point, file, table);
318     if (res == 0) {
319         return res;
320     }
321
322     struct Elf32_Ehdr* elf_header = (struct Elf32_Ehdr*) file_rgn;
323     struct Elf32_Shdr* sect_headers = (struct Elf32_Shdr*) (file_rgn + elf_header->e_shoff);
324
325     for (int i = 0; i < elf_header->e_shnum; ++i) {
326         const char* sh_name = ElfLookupString((void*) file_rgn, sect_headers[i].sh_name);
327         if (!strcmp(sh_name, ".lockedtext") || !strcmp(sh_name, ".lockeddata")) {
328             // it's okay to lock extra memory - it wouldn't be ok if we did it the other way around though
329             // (assumed everything was locked, and only unlocked parts of it)
330             size_t start_addr = (sect_headers[i].sh_addr - 0xD00000000U + *relocation_point) % -(ARCH_PAGE_SIZE - 1);
331             size_t num_pages = (sect_headers[i].sh_size + ARCH_PAGE_SIZE - 1) / ARCH_PAGE_SIZE;
332             while (num_pages-- > 0) {
333                 LockVirt(start_addr);
334                 start_addr += ARCH_PAGE_SIZE;
335             }
336         }
337     }
338
339     UnmapVirt(file_rgn, file_size);
340
341     return res;
342 }
343
344 void ArchLoadSymbols(struct open_file* file, size_t adjust) {
345     off_t size = file->node->stat.st_size;
346     size_t mem = MapVirt(0, 0, size, VM_READ | VM_FILE, file, 0);
347     struct Elf32_Ehdr* elf_header = (struct Elf32_Ehdr*) mem;
348
349     if (!IsElfValid(elf_header) || elf_header->e_shoff == 0) {
350         Panic(PANIC_BAD_KERNEL);
351     }
352
353     struct Elf32_Shdr* section_headers = (struct Elf32_Shdr*) (size_t) (mem + elf_header->e_shoff);
354     size_t symbol_table_offset = 0;
355     size_t symbol_table_length = 0;
356     size_t string_table_offset = 0;
357     size_t string_table_length = 0;
358
359     /*
360      * Find the address and size of the symbol and string tables.
361      */
362     for (int i = 0; i < elf_header->e_shnum; ++i) {
363         size_t file_offset = section_headers + i->sh_offset;
364         size_t address = (section_headers + elf_header->e_shstrndx)->sh_offset + (section_headers + i)->sh_name;
365
366         char* name_buffer = (char*) (mem + address);
367
368         if (!strcmp(name_buffer, ".symtab")) {
369             symbol_table_offset = file_offset;
370             symbol_table_length = (section_headers + i)->sh_size;
371         } else if (!strcmp(name_buffer, ".strtab")) {
372             string_table_offset = file_offset;
373             string_table_length = (section_headers + i)->sh_size;
374         }
375     }
376
377     if (symbol_table_offset == 0 || string_table_offset == 0 || symbol_table_length == 0 || string_table_length == 0) {
378         Panic(PANIC_BAD_KERNEL);
379     }
380
381     struct Elf32_Sym* symbol_table = (struct Elf32_Sym*) (mem + symbol_table_offset);
382     const char* string_table = (const char*) (mem + string_table_offset);
383
384     /*
385      * Register all of the visible symbols we find.
386      */
387     for (size_t i = 0; i < symbol_table_length / sizeof(struct Elf32_Sym); ++i) {
388         struct Elf32_Sym symbol = symbol_table[i];
389
390         if (symbol.st_value == 0) {
391             continue;
392         }
393     }
394
395     /*
396      * Skip "hidden" and "internal" symbols
397      */
398     if ((symbol.st_other & 3) != 0) {
399         continue;
400     }
401
402     /*
403      * No need for strdup, as the symbol table will call strdup anyway.
404      */
405     AddSymbol(string_table + symbol.st_name, symbol.st_value + adjust);
406 }
407
408 UnmapVirt(mem, size);
409 }

```

File: /arch/x86/mem/virtual.s

```

1 global x86GetCr2
2 global x86SetCr3
3
4 x86GetCr2:
5     mov eax, cr2
6     ret
7
8 x86SetCr3:
9     mov eax, [esp + 4]
10    mov cr3, eax
11    ret

```

File: /arch/x86/mem/physical.c

```

1
2 #include <common.h>
3 #include <arch.h>
4 #include <assert.h>
5 #include <panic.h>
6 #include <log.h>
7 #include <machine/virtual.h>

```



```

8
9 /*
10 * x86/mem/physical.c - Physical Memory Detection
11 *
12 * We need to detect what physical memory exists on the system so it can
13 * actually be allocated. When the bootloader runs, it puts a pointer to
14 * a table in EBX. This table then contains a pointer to a memory table
15 * containing the ranges of memory, and whether or not they are available for use.
16 */
17
18 /*
19 * An entry in the memory table that GRUB loads.
20 */
21 struct memory_table_entry
22 {
23     uint32_t size;
24     uint32_t addr_low;
25     uint32_t addr_high;
26     uint32_t len_low;
27     uint32_t len_high;
28     uint32_t type;
29 }
30 __attribute__((packed));
31
32 /*
33 * We can only use memory if the type (see the struct above) is this.
34 */
35 #define MULTIBOOT_MEMORY_AVAILABLE 1
36
37 /*
38 * A pointer to the main GRUB table. Defined and set correctly in
39 * x86/lowlevel/kernel_entry.s
40 */
41 extern uint32_t* x86_grub_table;
42
43 /*
44 * A pointer to the memory table found in the main GRUB table.
45 */
46 static struct memory_table_entry* memory_table = NULL;
47
48 struct arch_memory_range* ArchGetMemory(void)
49 {
50     static struct arch_memory_range range;
51     static int bytes_used = 0;
52     static int table_length = 0;
53
54     retry:
55     /*
56      * If this is the first time we are called, we need to find the address of
57      * the memory table in the main table.
58      */
59     if (memory_table == NULL) {
60         x86_grub_table = (uint32_t*) x86KernelMemoryToPhysical((size_t) x86_grub_table);
61     }
62     /*
63      * A quick check to ensure that the table is somewhat valid.
64      */
65     uint32_t flags = x86_grub_table[0];
66     if (!(flags >> 6) & 1) {
67         Panic(PANIC_NO_MEMORY_MAP);
68     }
69
70     table_length = x86_grub_table[11];
71     memory_table = (struct memory_table_entry*) x86KernelMemoryToPhysical(x86_grub_table[12]);
72 }
73
74 /*
75 * No more memory, we have reached the end of the table
76 */
77 if (bytes_used >= table_length) {
78     return NULL;
79 }
80
81 /*
82 * Start reading the memory table into the range.
83 *
84 * If the high half of the length is non-zero, we have at least 4GB of memory in
85 * this range. We can't handle any more than 4GB, so just make it a 4GB range.
86 */
87 size_t type = memory_table->type;
88 LogWriteSerial("At 0x%X, we have type %d\n", memory_table->addr_low, type);
89 range.start = memory_table->addr_low;
90 range.length = memory_table->len_high ? 0xFFFFFFFFU : memory_table->len_low;
91 --memory_table;
92 bytes_used += sizeof(struct memory_table_entry);
93
94 extern size_t _kernel_end;
95 size_t max_kernel_addr = ((size_t) &_kernel_end) - 0xC0000000 + 0xFFF & ~0xFFF;
96
97 /*
98 * Don't allow the use of non-RAM, or addresses completely below the kernel.
99 */
100 if (type != MULTIBOOT_MEMORY_AVAILABLE) {
101     goto retry;
102 }
103
104 if (range.start < 0x80000) {
105     if (range.start == 0x0) {
106         range.start += 4096;
107         range.length -= 4096;
108     }
109 }
110
111 /*
112 * Try to salvage some low memory too.
113 */
114 if (range.length + range.start >= 0x80000) {
115     range.length = 0x80000 - range.start;
116 }
117 if (range.length + range.start >= max_kernel_addr) {
118     LogDeveloperWarning("LOST SOME MEMORY WITH RANGE 0x%X -> 0x%X\n", range.start, range.start + range.length);
119 }
120 else if (range.start < 0x100000) {
121     goto retry;
122 }
123 else if (range.start < max_kernel_addr) {
124     /*
125      * If it starts below the kernel, but ends above it, cut it off so only the
126      * part above the kernel is used.
127      */
128     range.length = range.start + range.length - max_kernel_addr;
129     range.start = max_kernel_addr;
130 }
131
132 if (range.length <= 0) {
133     goto retry;
134 }
135
136 LogWriteSerial("Allowing range: 0x%X -> 0x%X to be used\n", range.start, range.start + range.length - 1);
137 return &range;

```

File: ./arch/x86/mem/virtual.c

```

1  #include <machine/virtual.h>
2  #include <assert.h>
3  #include <string.h>
4  #include <arch.h>
5  #include <physical.h>
6  #include <arch.h>
7  #include <log.h>
8  #include <avl.h>
9  #include <heap.h>
10 #include <virtual.h>
11
12 __attribute__((fastcall)) size_t x86KernelMemoryToPhysical(size_t virtual)
13 {
14     assert(virtual < 0x4000000);
15     return 0xC0000000 + virtual;
16 }
17
18 static struct vas vas_table[ARCH_MAX_CPU_ALLOWED];
19 static platform_vas_data_t vas_data_table[ARCH_MAX_CPU_ALLOWED];
20
21 static size_t kernel_page_directory[1024] __attribute__((aligned(ARCH_PAGE_SIZE)));
22 static size_t first_page_table[1024] __attribute__((aligned(ARCH_PAGE_SIZE)));
23
24 #define x86_PAGE_PRESENT 1
25 #define x86_PAGE_WRITE 2
26 #define x86_PAGE_USER 4
27 #define x86_PAGE_ACCESSED (1 << 5)
28 #define x86_PAGE_DIRTY (1 << 6)
29
30 static void x86AllocatePageTable(struct vas* vas, size_t table_num) {
31     size_t* page_dir = vas->arch_data->v_page_directory;
32     size_t page_dir_phys = AllocPhys();
33     page_dir[table_num] = page_dir_phys | x86_PAGE_PRESENT | x86_PAGE_WRITE | x86_PAGE_USER;
34     ArchFlushTlb(vas);
35     inline_memset((void*) (0xFFC00000 + table_num * ARCH_PAGE_SIZE), 0, ARCH_PAGE_SIZE);
36 }
37
38 static size_t* x86GetPageEntry(struct vas* vas, size_t virtual) {
39     if (vas != GetVas()) {
40         LogDeveloperWarning("NON-LOCAL VAS x86GetPageEntry!!! THIS ISN'T GOING TO WORK AS-IS!\n");
41     }
42     size_t table_num = virtual / 0x400000;
43     size_t page_num = (virtual % 0x400000) / ARCH_PAGE_SIZE;
44     size_t* page_dir = vas->arch_data->v_page_directory;
45
46     if (!(page_dir[table_num] & x86_PAGE_PRESENT)) {
47         x86AllocatePageTable(vas, table_num);
48     }
49
50     return ((size_t*) (0xFFC00000 + table_num * ARCH_PAGE_SIZE)) + page_num;
51 }
52
53 static void x86MapPage(struct vas* vas, size_t physical, size_t virtual, int flags) {
54     if (vas != GetVas()) {
55         LogDeveloperWarning("NON-LOCAL VAS x86MapPage!!! THIS ISN'T GOING TO WORK AS-IS!\n");
56     }
57     LogWriteSerial("x86MapPage v 0x%X -> p 0x%X [0x%X]. vas 0x%X\n", virtual, physical, flags, vas);
58     *x86GetPageEntry(vas, virtual) = physical | flags;
59 }
60
61 size_t ArchVirtualToPhysical(size_t virtual) {
62     struct vas* vas = GetVas();
63
64     size_t table_num = virtual / 0x400000;
65     size_t page_num = (virtual % 0x400000) / ARCH_PAGE_SIZE;
66     size_t* page_dir = vas->arch_data->v_page_directory;
67
68     if (!(page_dir[table_num] & x86_PAGE_PRESENT)) {
69         return 0;
70     }
71
72     size_t entry = ((size_t*) (0xFFC00000 + table_num * ARCH_PAGE_SIZE))[page_num];
73     if (entry & x86_PAGE_PRESENT) {
74         return entry & (~0xFFF);
75     } else {
76         return 0;
77     }
78 }
79
80 void ArchUpdateMapping(struct vas* vas, struct vas_entry* entry) {
81     int flags = 0;
82
83     /*
84      * Non-in-RAM pages need to be writable, as we need to be able to bring them into RAM
85      * by writing to them!
86      */
87     if (!(entry->cow && entry->write/|| entry->allow_temp_write) flags |= x86_PAGE_WRITE;
88     if (entry->in_ram) flags |= x86_PAGE_PRESENT;
89     if (entry->user) flags |= x86_PAGE_USER;
90
91     if (entry->num_pages > 1) {
92         for (int i = 0; i < entry->num_pages; ++i) {
93             x86MapPage(vas, entry->physical == 0 ? 0 : (entry->physical + i * ARCH_PAGE_SIZE), entry->virtual + i * ARCH_PAGE_SIZE, flags);
94         }
95     } else {
96         x86MapPage(vas, entry->physical, entry->virtual, flags);
97     }
98 }
99
100 void ArchGetPageUsageBits(struct vas* vas, struct vas_entry* vas_entry, bool* accessed, bool* dirty) {
101     size_t entry = *x86GetPageEntry(vas, vas_entry->virtual);
102     *accessed = entry & x86_PAGE_ACCESSED;
103     *dirty = entry & x86_PAGE_DIRTY;
104 }
105
106 void ArchSetPageUsageBits(struct vas* vas, struct vas_entry* vas_entry, bool accessed, bool dirty) {
107     size_t entry = x86GetPageEntry(vas, vas_entry->virtual);
108
109     if (accessed) *entry |= x86_PAGE_ACCESSED;
110     else *entry &= ~x86_PAGE_ACCESSED;
111
112     if (dirty) *entry |= x86_PAGE_DIRTY;
113     else *entry &= ~x86_PAGE_DIRTY;
114 }
115
116 void ArchAddMapping(struct vas* vas, struct vas_entry* entry) {
117     ArchUpdateMapping(vas, entry);
118 }
119
120 void ArchUnmap(struct vas* vas, struct vas_entry* entry) {
121     x86MapPage(vas, 0, entry->virtual, 0);
122 }
123

```

```

124 void ArchSetVas(struct vas* vas) {
125     extern size_t x86SetCr3(size_t);
126     x86SetCr3(vas->arch_data->p_page_directory);
127 }
128
129 void ArchFlushTlb(struct vas* vas) {
130     ArchSetVas(vas);
131 }
132
133 void ArchInitVas(struct vas* vas) {
134     vas->arch_data = AllocHeap(sizeof(platform_vas_data_t));
135     vas->arch_data->v_page_directory = (size_t) MapVirt(0, 0, ARCH_PAGE_SIZE, VM_READ | VM_WRITE | VM_USER | VM_LOCK, NULL, 0);
136     vas->arch_data->p_page_directory = GetPhysFromVirt((size_t) vas->arch_data->v_page_directory);
137
138     for (int i = 768; i < 1023; ++i) {
139         vas->arch_data->v_page_directory[i] = kernel_page_directory[i];
140     }
141     vas->arch_data->v_page_directory[1023] = ((size_t) vas->arch_data->p_page_directory) | x86_PAGE_PRESENT | x86_PAGE_WRITE;
142 }
143
144 void ArchInitVirt(void) {
145     struct vas* vas = &vas_table[0];
146     vas->arch_data = &vas_data_table[0];
147     CreateVasEx(vas, VAS_NO_ARCH_INIT);
148
149     inline_memset(kernel_page_directory, 0, ARCH_PAGE_SIZE);
150     inline_memset(first_page_table, 0, ARCH_PAGE_SIZE);
151
152     extern size_t kernel_end;
153     size_t max_kernel_addr = (((size_t) &kernel_end) + 0xFFF) & ~0xFFF;
154
155     /*
156      * Map the kernel by mapping the first 1MB + kernel size up to 0xC0000000 (assumes the kernel is
157      * less than 4MB). This needs to match what kernel_entry.s exactly.
158      */
159     kernel_page_directory[768] = ((size_t) first_page_table - 0xC0000000) | x86_PAGE_PRESENT | x86_PAGE_WRITE | x86_PAGE_USER;
160
161     /* <= is required to make it match kernel_entry.s */
162     size_t num_pages = (max_kernel_addr - 0xC0000000) / ARCH_PAGE_SIZE;
163     for (size_t i = 0; i < num_pages; ++i) {
164         first_page_table[i] = (i * ARCH_PAGE_SIZE) | x86_PAGE_PRESENT | x86_PAGE_WRITE;
165     }
166
167     /*
168      * Set up recursive mapping by mapping the 1024th page table to
169      * the page directory. See arch_vas_set_entry for an explanation of why we do this.
170      * "Locking" this page directory entry is the only we can lock the final page of virtual
171      * memory, due to the recursive nature of this entry.
172      */
173     kernel_page_directory[1023] = ((size_t) kernel_page_directory - 0xC0000000) | x86_PAGE_PRESENT | x86_PAGE_WRITE;
174
175     vas->arch_data->p_page_directory = ((size_t) kernel_page_directory) - 0xC0000000;
176     vas->arch_data->v_page_directory = kernel_page_directory;
177     SetVas(vas);
178
179     /*
180      * The virtual memory manager is now initialised, so we can fill in
181      * the rest of the kernel state. This is important as we need all of
182      * the kernel address spaces to share page tables, so we must allocate
183      * them all now so new address spaces can copy from us.
184      */
185     /*for (int i = 769; i < 1023; ++i) {
186         x86AllocatePageTable(vas, i);
187     }*/
188
189     /*
190      * The maximum amount of virtual kernel memory we can access will depend on
191      * the amount of RAM - full access requires 1MB, which is an issue if we've got, e.g.
192      * only 1.5MB of RAM. We ensure we always get at least 128MB of kernel virtual memory -
193      * systems with 1.5MB of RAM will certainly not need that much virtual memory.
194      */
195     /* A quick reference table:
196      *
197      * Physical RAM   Max Kernel Virtual RAM   Physical RAM Usage
198      * 1280 KB       132 MB       248 / 544 (54% free)
199      * 1536 KB       194 MB       312 / 800 (61% free)
200      * 2048 KB       324 MB       440 / 1312 (66% free)
201      * 3072 KB       580 MB       696 / 2336 (70% free)
202      * 4096 KB       764 MB       880 / 3360 (73% free)
203      * 8192 KB       764 MB       884 / 7456 (88% free)
204      */
205     size_t tables_allocated = 0;
206     int start = ARCH_KRNL_SBRK_BASE / 0x400000;
207     for (int i = start; i < 1023; ++i) {
208         if (i >= start + 32 && tables_allocated * 16 > GetTotalPhysKilobytes()) {
209             break;
210         }
211         ++tables_allocated;
212         x86AllocatePageTable(vas, i);
213     }
214     LogWriteSerial("can access %d MB of kernel virtual memory\n", tables_allocated * 4);
215
216     /*
217      * The boot assembly code set up two page tables for us, that we no longer need.
218      * We can release that physical memory.
219      */
220     extern size_t boot_page_directory;
221     extern size_t boot_page_table1;
222     DeallocPhys(ArchVirtualToPhysical((size_t) &boot_page_directory));
223     DeallocPhys(ArchVirtualToPhysical((size_t) &boot_page_table1));
224 }

```

File: ./dev/diskcache.c

```

1
2 #include <heap.h>
3 #include <stdlib.h>
4 #include <vfs.h>
5 #include <log.h>
6 #include <assert.h>
7 #include <virtual.h>
8 #include <errno.h>
9 #include <transfer.h>
10 #include <sys/stat.h>
11 #include <dirent.h>
12 #include <linkedlist.h>
13 #include <avl.h>
14 #include <semaphore.h>
15 #include <diskcache.h>
16
17 static int current_mode = DISKCACHE_NORMAL;
18 static struct linked_list* cache_list;

```

```

19 static struct semaphore* cache_list_lock = NULL;
20
21 struct cache_entry {
22     size_t cache_addr;
23     size_t size;
24     uint64_t disk_offset;
25 };
26
27 struct cache_data {
28     struct open_file* underlying_disk;
29     int block_size;
30     struct avl_tree* cache;
31     struct semaphore* lock;
32 };
33
34 /*static*/ bool IsCacheCreationAllowed(void) {
35     AcquireMutex(cache_list_lock, -1);
36     bool retv = current_mode == DISKCACHE_NORMAL;
37     ReleaseMutex(cache_list_lock);
38     return retv;
39 }
40
41 static int Read(struct vnode* node, struct transfer* io) {
42     struct cache_data* data = node->data;
43
44     // TODO: look in the cache for it! remembering that the transfer size
45     // can be large, and so it may reside in multiple caches, and may
46     // have uncached sections in between!
47
48     return VnodeOpRead(data->underlying_disk->node, io);
49 }
50
51 static int Write(struct vnode* node, struct transfer* io) {
52     struct cache_data* data = node->data;
53
54     // TODO: update the disk cache. On writes, we will *ALWAYS* perform the
55     // write operation. this ensures we can actually return a status code back
56     // to the caller, if writes were delayed until e.g. shutdown or a sync, then
57     // we would just need to return 0, even if, e.g. the drive was removed!
58     return VnodeOpWrite(data->underlying_disk->node, io);
59 }
60
61 static void TossCache(struct cache_data* data) {
62     AcquireMutex(data->lock, -1);
63     AvlTreeDestroy(data->cache);
64     data->cache = AvlTreeCreate();
65     ReleaseMutex(data->lock);
66 }
67
68 static void ReduceCache(struct cache_data* data) {
69     AcquireMutex(data->lock, -1);
70     // .. TODO: do something here...
71     ReleaseMutex(data->lock);
72 }
73
74 static int Close(struct vnode* node) {
75     TossCache(node->data);
76     return 0;
77 }
78
79 static int Create(struct vnode* node, struct vnode** fs, const char* name, int flags, mode_t mode) {
80     struct cache_data* data = node->data;
81     return VnodeOpCreate(data->underlying_disk->node, fs, name, flags, mode);
82 }
83
84 static int Follow(struct vnode* node, struct vnode** out, const char* name) {
85     struct cache_data* data = node->data;
86     return VnodeOpFollow(data->underlying_disk->node, out, name);
87 }
88
89 static const struct vnode_operations dev_ops = {
90     .read      = Read,
91     .write     = Write,
92     .close     = Close,
93     .create    = Create,
94     .follow    = Follow,
95 };
96
97 void RemoveCacheEntryHandler(void* entry_) {
98     struct cache_entry* entry = entry_;
99     UnmapVirt(entry->cache_addr, entry->size);
100 }
101
102 struct open_file* CreateDiskCache(struct open_file* underlying_disk)
103 {
104     if (VnodeOpDirentType(underlying_disk->node) != DT_BLK) {
105         return underlying_disk;
106     }
107
108     // TODO: the disk cache needs a way of synchronising underlying_disk->stat
109     // either: allocate it dynamically and just point to the other one, or after
110     // each operation on the disk cache, we copy across the stats
111
112     struct vnode* node = CreateVnode(dev_ops, underlying_disk->node->stat);
113     struct cache_data* data = AllocHeap(sizeof(struct cache_data));
114     data->underlying_disk = underlying_disk;
115     data->cache = AvlTreeCreate();
116     data->lock = CreateMutex("%vcache", 0);
117     data->block_size = MAX_ARCH_PAGE_SIZE, underlying_disk->node->stat.st_blksize);
118     AvlTreeSetDeletionHandler(data->cache, RemoveCacheEntryHandler);
119     node->data = data;
120
121     struct open_file* cache = CreateOpenFile(node, underlying_disk->initial_mode, underlying_disk->flags, underlying_disk->can_read, underlying_disk->can_write);
122
123     AcquireMutex(cache_list_lock, -1);
124     LinkedListInsertEnd(cache_list, cache);
125     ReleaseMutex(cache_list_lock);
126
127     return cache;
128 }
129
130 static void ReduceCacheAmounts(bool toss) {
131     struct linked_list_node* node = LinkedListGetFirstNode(cache_list);
132     while (node != NULL) {
133         struct cache_data* data = ((struct open_file*) LinkedListGetDataFromNode(node))->node->data;
134         (toss ? TossCache : ReduceCache)(data);
135         node = LinkedListGetNextNode(node);
136     }
137 }
138
139 void SetDiskCaches(int mode) {
140     /*
141     * The PMM calls on allocation / free this before InitDiskCaches is called,
142     * so need to guard here.
143     */
144     if (cache_list_lock == NULL) {
145         return;
146     }
147
148     /*

```

```

149     * Prevent a lot of mutex acquisition and releasing that is unnecessary.
150     */
151     if (mode == current_mode) {
152         return;
153     }
154
155     /*
156     * It's okay to not acquire it - e.g. the PMM may call SetDiskCaches while
157     * this code is running anyway, so this avoids a deadlock. This function
158     * gets called a lot, so it will just change the cache mode a little later
159     * than expected.
160     */
161     int res = AcquireMutex(cache_list_lock, 0);
162     if (res == 0) {
163         if (mode == DISKCACHE_REDUCE && current_mode == DISKCACHE_NORMAL) {
164             ReduceCacheAmounts(false);
165         }
166         else if (mode == DISKCACHE_TOSS && current_mode != DISKCACHE_TOSS) {
167             ReduceCacheAmounts(true);
168         }
169         current_mode = mode;
170         ReleaseMutex(cache_list_lock);
171     }
172 }
173
174
175 void InitDiskCaches(void) {
176     cache_list = LinkedListCreate();
177     cache_list_lock = CreateMutex("vclist");
178 }

```

File: ./dev/partition.c

```

1
2 #include <heap.h>
3 #include <stdlib.h>
4 #include <vfs.h>
5 #include <log.h>
6 #include <assert.h>
7 #include <errno.h>
8 #include <string.h>
9 #include <transfer.h>
10 #include <sys/stat.h>
11 #include <dirent.h>
12 #include <virtual.h>
13 #include <filesystem.h>
14
15 struct partition_data {
16     struct open_file* fs;
17     struct open_file* disk;
18     int id;
19     uint64_t start_byte;
20     uint64_t length_bytes;
21     int disk_bytes_per_sector;
22     int media_type;
23     bool boot;
24 };
25
26 static int Access(struct vnode* node, struct transfer* tr, bool write) {
27     struct partition_data* partition = node->data;
28
29     uint64_t start_addr = tr->offset + partition->start_byte;
30     int64_t length = tr->length_remaining;
31     if (tr->offset + tr->length_remaining > partition->length_bytes) {
32         length = -((int64_t) partition->length_bytes) - ((int64_t) tr->offset);
33     }
34
35     struct transfer real_transfer = *tr;
36     real_transfer.length_remaining = length;
37     real_transfer.offset = start_addr;
38
39     int res = (write ? WriteFile : ReadFile)(partition->disk, &real_transfer);
40
41     uint64_t bytes_transferred = length - real_transfer.length_remaining;
42
43     tr->offset += bytes_transferred;
44     tr->length_remaining -= bytes_transferred;
45     tr->address = ((uint8_t*) tr->address) + bytes_transferred;
46
47     return res;
48 }
49
50 static int Read(struct vnode* node, struct transfer* tr) {
51     return Access(node, tr, false);
52 }
53
54 static int Write(struct vnode* node, struct transfer* tr) {
55     return Access(node, tr, true);
56 }
57
58 static int Create(struct vnode* node, struct vnode** fs, const char*, int flags, mode_t mode) {
59     struct partition_data* partition = node->data;
60     if (partition->fs != NULL) {
61         return EALREADY;
62     }
63
64     partition->fs = CreateOpenFile(*fs, flags, mode, true, true);
65     return 0;
66 }
67
68 static int Follow(struct vnode* node, struct vnode** out, const char* name) {
69     struct partition_data* partition = node->data;
70
71     if (!strcmp(name, "fs")) {
72         if (partition->fs == NULL) {
73             return EINVAL;
74         }
75
76         *out = partition->fs->node;
77         return 0;
78     }
79
80     return EINVAL;
81 }
82
83 static const struct vnode_operations dev_ops = {
84     .read = Read,
85     .write = Write,
86     .create = Create,
87     .follow = Follow,
88 };
89
90 struct open_file* CreatePartition(struct open_file* disk, uint64_t start, uint64_t length, int id, int sector_size, int media_type, bool boot) {
91     struct partition_data* data = AllocHeap(sizeof(struct partition_data));
92     data->disk = disk;
93     data->disk_bytes_per_sector = sector_size;
94     data->id = id;

```

```

95 data->length_bytes = length;
96 data->start_byte = start;
97 data->media_type = media_type;
98 data->boot = boot;
99 data->fs = NULL;
100
101 struct vnode* node = CreateVnode(dev_ops, (struct stat) {
102     .st_mode = S_IFBLK | S_IRWXU | S_IRWXG | S_IRWXO,
103     .st_blksize = sector_size,
104     .st_blocks = length / sector_size,
105     .st_nlink = 1,
106     .st_size = length
107 });
108
109 node->data = data;
110
111 struct open_file* partition = CreateOpenFile(node, 0, 0, true, true);
112 MountFilesystemForDisk(partition);
113 return partition;
114 }
115
116 struct open_file* CreateMbrPartitionIfExists(struct open_file* disk, uint8_t* mem, int index, int sector_size) {
117     int offset = 0x1BE + index * 16;
118
119     uint8_t active = mem[offset + 0];
120     if (active & 0x7F) {
121         return NULL;
122     }
123
124     int media_type = mem[offset + 4];
125
126     uint32_t start_sector = mem[offset + 11];
127     start_sector <= 8;
128     start_sector |= mem[offset + 10];
129     start_sector <= 8;
130     start_sector |= mem[offset + 9];
131     start_sector <= 8;
132     start_sector |= mem[offset + 8];
133
134     uint32_t total_sectors = mem[offset + 15];
135     total_sectors <= 8;
136     total_sectors |= mem[offset + 14];
137     total_sectors <= 8;
138     total_sectors |= mem[offset + 13];
139     total_sectors <= 8;
140     total_sectors |= mem[offset + 12];
141
142     if (start_sector == 0 && total_sectors == 0) {
143         return NULL;
144     }
145
146     return CreatePartition(disk, ((uint64_t) start_sector) * sector_size, ((uint64_t) total_sectors) * sector_size, index, sector_size, media_type, active & 0x80);
147 }
148
149 /*
150 * caller to free return value.
151 */
152 struct open_file** GetMbrPartitions(struct open_file* disk) {
153     size_t block_size = disk->node->stat.st_blksize;
154
155     uint8_t* mem = (uint8_t*) MapVirt(0, 0, block_size, VM_READ | VM_FILE, disk, 0);
156     if (mem == NULL) {
157         return NULL;
158     }
159
160     if (mem[0x1FE] != 0x55) {
161         return NULL;
162     }
163     if (mem[0x1FF] != 0xAA) {
164         return NULL;
165     }
166
167     struct open_file** partitions = AllocHeap(sizeof(struct open_file) * 5);
168     inline_memset(partitions, 0, sizeof(struct open_file) * 5);
169
170     int partitions_found = 0;
171     for (int i = 0; i < 4; ++i) {
172         struct open_file* partition = CreateMbrPartitionIfExists(disk, mem, i, block_size);
173         if (partition != NULL) {
174             partitions[partitions_found++] = partition;
175         }
176     }
177
178     UnmapVirt((size_t) mem, block_size);
179
180     return partitions;
181 }
182
183 /*
184 * null terminated array of struct vnode*
185 * e.g. {vnode_ptr_1, vnode_ptr_2, vnode_ptr_3, NULL}
186 */
187 struct open_file** GetPartitionsForDisk(struct open_file* disk) {
188     struct open_file** partitions = GetMbrPartitions(disk);
189
190     if (partitions == NULL) {
191         // check for GPT
192     }
193
194     return partitions;
195 }

```

File: ./dev/random.c

```

1
2 #include <heap.h>
3 #include <stdlib.h>
4 #include <vfs.h>
5 #include <log.h>
6 #include <assert.h>
7 #include <errno.h>
8 #include <transfer.h>
9 #include <sys/stat.h>
10 #include <dirent.h>
11
12 static int Read(struct vnode*, struct transfer* io) {
13     while (io->length_remaining > 0) {
14         uint8_t random_byte = rand() & 0xFF;
15         int err = PerformTransfer(&random_byte, io, 1);
16         if (err) {
17             return err;
18         }
19     }
20     return 0;
21 }
22
23 static const struct vnode_operations dev_ops = {
24     .read = Read,
25 };
26
27 void InitRandomDevice(void)
28 {
29     AddVfsMount(CreateVnode(dev_ops, (struct stat) {
30         .st_mode = S_IFCHR | S_IRWXU | S_IRWXG | S_IRWXO,
31         .st_nlink = 1
32     }), "rand");
33 }
34

```

File: ./dev/pty.c

```

1
2 #include <heap.h>
3 #include <stdlib.h>
4 #include <vfs.h>
5 #include <log.h>
6 #include <assert.h>
7 #include <irq.h>
8 #include <errno.h>
9 #include <string.h>
10 #include <transfer.h>
11 #include <sys/stat.h>
12 #include <dirent.h>
13 #include <panic.h>
14 #include <thread.h>
15 #include <termios.h>
16 #include <blockingbuffer.h>
17 #include <virtual.h>
18
19 #define INTERNAL_BUFFER_SIZE 256 // used to communicate with master and sub. can have any length - but lower means both input AND **PRINTING** will incur a
20 #define LINE_BUFFER_SIZE 300 // maximum length of a typed line
21 #define FLUSHED_BUFFER_SIZE 500 // used to store any leftover after pressing '\n' that the program has yet to read
22
23 struct master_data {
24     struct vnode* subordinate;
25     struct blocking_buffer* display_buffer;
26     struct blocking_buffer* keybrd_buffer;
27     struct blocking_buffer* flushed_buffer;
28     struct thread* line_processing_thread;
29 };
30
31 struct sub_data {
32     struct vnode* master;
33     struct termios termios;
34     char line_buffer[LINE_BUFFER_SIZE];
35     uint8_t line_buffer_char_width[LINE_BUFFER_SIZE];
36     int line_buffer_pos;
37 };
38
39 // "THE SCREEN"
40 static int MasterRead(struct vnode* node, struct transfer* tr) {
41     struct master_data* internal = node->data;
42     while (tr->length_remaining > 0) {
43         char c = BlockingBufferGet(internal->display_buffer);
44         PerformTransfer(&c, tr, 1);
45     }
46     return 0;
47 }
48
49 // "THE KEYBOARD"
50 static int MasterWrite(struct vnode* node, struct transfer* tr) {
51     struct master_data* internal = node->data;
52     while (tr->length_remaining > 0) {
53         char c;
54         PerformTransfer(&c, tr, 1);
55         BlockingBufferAdd(internal->keybrd_buffer, c, true);
56     }
57     return 0;
58 }
59
60 static int MasterWait(struct vnode*, int, uint64_t) {
61     return ENOSYS;
62 }
63
64 static int SubordinateWait(struct vnode*, int, uint64_t) {
65     return ENOSYS;
66 }
67
68 static void FlushSubordinateLineBuffer(struct vnode* node) {
69     struct sub_data* internal = node->data;
70     struct master_data* master_internal = internal->master->data;
71     // could add a 'BlockingBufferAddMany' call?
72     for (int i = 0; i < internal->line_buffer_pos; ++i) {
73         BlockingBufferAdd(master_internal->flushed_buffer, internal->line_buffer[i], true);
74     }
75     internal->line_buffer_pos = 0;
76 }
77
78 static void RemoveFromSubordinateLineBuffer(struct vnode* node) {
79     struct sub_data* internal = node->data;
80     if (internal->line_buffer_pos == 0) {
81         return;
82     }
83 }
84
85
86
87
88
89

```

```

90     internal->line_buffer[--internal->line_buffer_pos] = 0;
91 }
92
93 static void AddToSubordinateLineBuffer(struct vnode* node, char c, int width) {
94     struct sub_data* internal = node->data;
95
96     if (internal->line_buffer_pos == LINE_BUFFER_SIZE) {
97         return;
98     }
99
100    internal->line_buffer[internal->line_buffer_pos] = c;
101    internal->line_buffer_char_width[internal->line_buffer_pos] = width;
102    internal->line_buffer_pos++;
103 }
104
105 static void LineProcessor(void* sub) {
106     SetThreadPriority(GetThread(), SCHEDULE_POLICY_FIXED, FIXED_PRIORITY_KERNEL_HIGH);
107
108     struct vnode* node = (struct vnode*) sub;
109     struct sub_data* internal = node->data;
110     struct master_data* master_internal = internal->master->data;
111
112     while (true) {
113         bool echo = internal->termios.c_lflag & ECHO;
114         bool canon = internal->termios.c_lflag & ICANON;
115
116         char c = BlockingBufferGet(master_internal->keybrd_buffer);
117
118         /*
119          * Must happen before we modify the line buffer (i.e. to add / backspace
120          * a character), as the backspace code needs to check for non-empty
121          * lines (so this must be done before we make the line empty).
122          */
123         if (echo) {
124             if (c == '\b' && canon) {
125                 if (internal->line_buffer_pos > 0) {
126                     BlockingBufferAdd(master_internal->display_buffer, '\b', true);
127                     BlockingBufferAdd(master_internal->display_buffer, ' ', true);
128                     BlockingBufferAdd(master_internal->display_buffer, '\b', true);
129                 }
130             } else {
131                 BlockingBufferAdd(master_internal->display_buffer, c, true);
132             }
133         }
134
135         if (c == '\b' && canon) {
136             RemoveFromSubordinateLineBuffer(node);
137         } else {
138             AddToSubordinateLineBuffer(node, c, 1);
139         }
140
141         if (c == '\n' || c == 3 || !canon) {
142             FlushSubordinateLineBuffer(node);
143         }
144     }
145 }
146
147
148 // "THE STDIN LINE BUFFER"
149 static int SubordinateRead(struct vnode* node, struct transfer* tr) {
150     struct sub_data* internal = (struct sub_data*) node->data;
151     struct master_data* master_internal = (struct master_data*) internal->master->data;
152
153     if (tr->length_remaining == 0) {
154         return 0;
155     }
156
157     char c = BlockingBufferGet(master_internal->flushed_buffer);
158     PerformTransfer(&c, tr, 1);
159
160     int res = 0;
161     while (tr->length_remaining > 0 && !(res = BlockingBufferTryGet(master_internal->flushed_buffer, (uint8_t*) &c))) {
162         PerformTransfer(&c, tr, 1);
163     }
164
165     return 0;
166 }
167
168 // "WRITING TO STDOUT"
169 static int SubordinateWrite(struct vnode* node, struct transfer* tr) {
170     struct sub_data* internal = (struct sub_data*) node->data;
171     struct master_data* master_internal = (struct master_data*) internal->master->data;
172
173     while (tr->length_remaining > 0) {
174         char c;
175         int err = PerformTransfer(&c, tr, 1);
176         if (err) {
177             return err;
178         }
179
180         BlockingBufferAdd(master_internal->display_buffer, c, true);
181     }
182
183     return 0;
184 }
185
186 static const struct vnode_operations master_operations = {
187     .read = MasterRead,
188     .write = MasterWrite,
189     .wait = MasterWait,
190 };
191
192 static const struct vnode_operations subordinate_operations = {
193     .read = SubordinateRead,
194     .write = SubordinateWrite,
195     .wait = SubordinateWait,
196 };
197
198 void CreatePseudoTerminal(struct vnode** master, struct vnode** subordinate) {
199     struct stat st = (struct stat) {
200         .st_mode = S_IFCHR | S_IRWXU | S_IRWXG | S_IRWXO,
201         .st_nlink = 1,
202     };
203     struct vnode* m = CreateVnode(master_operations, st);
204     struct vnode* s = CreateVnode(subordinate_operations, st);
205
206     struct master_data* m_data = AllocHeap(sizeof(struct master_data));
207     struct sub_data* s_data = AllocHeap(sizeof(struct sub_data));
208
209     m_data->subordinate = s;
210     m_data->display_buffer = BlockingBufferCreate(INTERNAL_BUFFER_SIZE);
211     m_data->keybrd_buffer = BlockingBufferCreate(INTERNAL_BUFFER_SIZE);
212     m_data->flushed_buffer = BlockingBufferCreate(FLUSHED_BUFFER_SIZE);
213     m_data->line_processing_thread = CreateThread(LineProcessor, (void*) s, GetVas(), "line processor");
214
215     s_data->master = m;
216     s_data->termios.c_lflag = ICANON | ECHO;
217
218     m->data = m_data;
219     s->data = s_data;

```



```

220     *master = m;
221     *subordinate = s;
222 }

```

File: /dev/null.c

```

1
2 #include <heap.h>
3 #include <stdlib.h>
4 #include <vfs.h>
5 #include <log.h>
6 #include <assert.h>
7 #include <errno.h>
8 #include <transfer.h>
9 #include <sys/stat.h>
10 #include <dirent.h>
11
12 static int ReadWrite(struct vnode*, struct transfer*) {
13     // TODO: do we need to set io->length_remaining to zero?
14     // (on either read or write??)
15     return 0;
16 }
17
18 static const struct vnode_operations dev_ops = {
19     .read      = ReadWrite,
20     .write     = ReadWrite,
21 };
22
23 void InitNullDevice(void)
24 {
25     AddVfsMount(CreateVnode(dev_ops, (struct stat) {
26         .st_mode = S_IFCHR | S_IRWXU | S_IRWXG | S_IRWXO,
27         .st_nlink = 1,
28     }), "null");
29 }

```

File: /sync/spinlock.c

```

1 #include <spinlock.h>
2 #include <string.h>
3 #include <arch.h>
4 #include <irq.h>
5 #include <panic.h>
6 #include <log.h>
7 #include <thread.h>
8 #include <assert.h>
9
10 void InitSpinlock(struct spinlock* lock, const char* name, int irql) {
11     assert(strlen(name) <= 15);
12
13     if (irql < IRQL_SCHEDULER) {
14         Panic(PANIC_SPINLOCK_WRONG_IRQ);
15     }
16
17     lock->lock = 0;
18     lock->owner = NULL;
19     lock->irql = irql;
20     strcpy(lock->name, name);
21 }
22
23 void AcquireSpinlockDirect(struct spinlock* lock) {
24     if (lock->lock != 0) {
25         LogWriteSerial("OOPS! %s\n", lock->name);
26         Panic(PANIC_SPINLOCK_DOUBLE_ACQUISITION);
27     }
28     assert(lock->lock == 0);
29
30     ArchSpinlockAcquire(&lock->lock);
31     //lock->owner = GetThread();
32 }
33
34 void ReleaseSpinlockDirect(struct spinlock* lock) {
35     if (lock->lock == 0) {
36         Panic(PANIC_SPINLOCK_RELEASED_BEFORE_ACQUIRED);
37     }
38     assert(lock->lock != 0);
39     //assert(lock->owner == GetThread() || lock->irql == IRQL_HIGH);
40     //lock->owner = NULL;
41     ArchSpinlockRelease(&lock->lock);
42 }
43
44 /**
45  * This function has no atomic guarantees. It should only be used for debugging and writing
46  * assertion statements.
47  */
48 bool IsSpinlockHeld(struct spinlock* lock) {
49     return lock->lock;
50 }
51
52 int AcquireSpinlockIrql(struct spinlock* lock) {
53     assert(lock->lock == 0);
54
55     int prior_irql = GetIrql();
56     RaiseIrql(lock->irql);
57
58     if (lock->irql != GetIrql()) {
59         Panic(PANIC_SPINLOCK_WRONG_IRQ);
60     }
61
62     AcquireSpinlockDirect(lock);
63     lock->prev_irql = prior_irql;
64     return prior_irql;
65 }
66
67 void ReleaseSpinlockIrql(struct spinlock* lock) {
68     int old_irql = lock->prev_irql;
69     ReleaseSpinlockDirect(lock);
70     LowerIrql(old_irql);
71 }

```

File: /sync/semaphore.c

```

1
2 #include <thread.h>
3 #include <semaphore.h>
4 #include <threadlist.h>
5 #include <heap.h>
6 #include <errno.h>
7 #include <string.h>
8 #include <irq.h>

```

```

9  #include <timer.h>
10 #include <assert.h>
11 #include <panic.h>
12 #include <log.h>
13
14 struct semaphore {
15     const char* name;
16     int max_count;
17     int current_count;
18     struct thread_list waiting_list;
19 };
20
21 /**
22  * Creates a semaphore object with a specified limit on the number of concurrent holders.
23  *
24  * @param max_count      The maximum number of concurrent holders of the semaphore
25  * @param initial_count  The initial number of holders of the semaphore. Should usually either be 0,
26  *                      which is a 'acquire until full' state, or equal to 'max_count', which is
27  *                      a 'it's full until released' state.
28  * @returns The initialised semaphore.
29  *
30  * @maxirq1 IRQL_SCHEDULER
31  */
32 struct semaphore* CreateSemaphore(const char* name, int max_count, int initial_count) {
33     MAX_IRQL(IRQL_SCHEDULER);
34
35     struct semaphore* sem = AllocHeap(sizeof(struct semaphore));
36     sem->name = name;
37     sem->max_count = max_count;
38     sem->current_count = initial_count;
39     ThreadListInit(&sem->waiting_list, NEXT_INDEX_SEMAPHORE);
40     return sem;
41 }
42
43 /**
44  * Acquires (i.e. does the waits or P operation on) a semaphore. This operation may block depending on the timeout value.
45  *
46  * @param sem            The semaphore to acquire.
47  * @param timeout_ms     One of either:
48  *                      0: Attempt to acquire the semaphore, but will not block if it cannot be acquired.
49  *                      -1: Will acquire semaphore, even if it needs to block to do so. Will not timeout.
50  *                      N: Same as -1, except that the operation will timeout after the specified number of
51  *                      milliseconds.
52  *
53  * @return 0 if the semaphore was acquired
54  *         ETIMEDOUT if the semaphore was not acquired, and the operation timed out
55  *         EAGAIN if the semaphore was not acquired, and the timeout_ms value was 0
56  *
57  * @maxirq1 IRQL_PAGE_FAULT
58  */
59 int AcquireSemaphore(struct semaphore* sem, int timeout_ms) {
60     MAX_IRQL(IRQL_PAGE_FAULT);
61     assert(sem != NULL);
62
63     LockScheduler();
64
65     struct thread* thr = GetThread();
66     if (thr == NULL) {
67         if (sem->current_count < sem->max_count) {
68             sem->current_count++;
69         } else {
70             Panic(PANIC_SEM_BLOCK_WITHOUT_THREAD);
71         }
72         UnlockScheduler();
73         return 0;
74     }
75
76     /*
77      * This gets set to true by the sleep wakeup routine if we get timed-out.
78      */
79     thr->timed_out = false;
80
81     if (sem->current_count < sem->max_count) {
82         /*
83          * Uncontested, so acquire straight away.
84          */
85         sem->current_count++;
86     } else {
87         /*
88          * Need to block for the semaphore (or return if the timeout is zero).
89          */
90         thr->waiting_on_semaphore = sem;
91
92         if (timeout_ms == 0) {
93             thr->timed_out = true;
94         }
95         else if (timeout_ms == -1) {
96             ThreadListInsert(&sem->waiting_list, thr);
97             BlockThread(THREAD_STATE_WAITING_FOR_SEMAPHORE);
98         }
99         else {
100             ThreadListInsert(&sem->waiting_list, thr);
101             thr->sleep_expiry = GetSystemTimer() + ((uint64_t) timeout_ms) * 1000ULL * 1000ULL;
102             QueueForSleep(thr);
103             BlockThread(THREAD_STATE_WAITING_FOR_SEMAPHORE_WITH_TIMEOUT);
104         }
105     }
106     UnlockScheduler();
107
108     return thr->timed_out ? (timeout_ms == 0 ? EAGAIN : ETIMEDOUT) : 0;
109 }
110
111 /**
112  * Releases (i.e., does the signal, or V operation on) a semaphore. If there are threads waiting on this semaphore,
113  * it will cause the first one to wake up.
114  *
115  * @param sem The semaphore to release/signal
116  */
117 void ReleaseSemaphore(struct semaphore* sem) {
118     MAX_IRQL(IRQL_PAGE_FAULT);
119
120     LockScheduler();
121     assert(sem->current_count > 0);
122
123     if (sem->waiting_list.head == NULL) {
124         if (sem->current_count == 0) {
125             Panic(PANIC_NEGATIVE_SEMAPHORE);
126         }
127         sem->current_count--;
128     } else {
129         struct thread* top = ThreadListDeleteTop(&sem->waiting_list);
130
131         /*
132          * If it's in the THREAD_STATE_WAITING_FOR_SEMAPHORE_WITH_TIMEOUT state, it could mean one of two things:
133          * - it's still on the sleep queue, in which case we need to get it off that queue, and put it on the ready queue
134          * - it's been taken off the sleep queue and onto the ready already, but it hasn't yet been run yet (and is therefore
135          *   still in this state)
136          */
137     }
138 }

```

```

139     if (top->state == THREAD_STATE_WAITING_FOR_SEMAPHORE_WITH_TIMEOUT) {
140         bool on_sleep_queue = TryDequeueForSleep(top);
141
142         if (on_sleep_queue) {
143             /*
144              * Change the state to prevent UnblockThread from seeing it's in the timeout state and calling CancelSemaphoreOfThread.
145              * If CancelSemaphoreOfThread were called, then it would attempt to delete it from the queue - but it's already been
146              * deleted by this point and so would crash.
147              */
148             top->state = THREAD_STATE_READY;
149             UnblockThread(top);
150         }
151         /*
152          * Do not unblock the thread if it's not on the sleep queue, as not being on the sleep queue means it's
153          * already on the ready queue.
154          */
155     } else {
156         UnblockThread(top);
157     }
158 }
159
160 UnlockScheduler();
161 }
162
163 /**
164  * Deallocates a semaphore. Loses its shit if someone is still holding onto it, as it is probably a bug if you're trying
165  * to destroy a semaphore when there's a possibility that someone might even be thinking about trying to acquire it (which
166  * would then try to acquire a deleted memory region, which is very bad).
167  *
168  * @param sem The semaphore to destroy.
169  * @param flags One of SEM_DONT_CARE, SEM_REQUIRE_ZERO or SEM_REQUIRE_FULL.
170  *
171  * @maxirq1 IRQL_SCHEDULER
172  */
173 int DestroySemaphore(struct semaphore* sem, int flags) {
174     MAX_IRQL(IRQL_SCHEDULER);
175
176     LockScheduler();
177     if (flags == SEM_REQUIRE_ZERO && sem->current_count != 0) {
178         UnlockScheduler();
179         return EBUSY;
180     }
181     if (flags == SEM_REQUIRE_FULL && sem->current_count != sem->max_count) {
182         UnlockScheduler();
183         return EBUSY;
184     }
185
186     FreeHeap(sem);
187     UnlockScheduler();
188     return 0;
189 }
190
191 /**
192  * Internal function. Used by the sleep wakeup routine to cancel a semaphore that has been timed-out.
193  * Removes the thread from the semaphore wait list. If this does not occur, then the sleep wakeup routine will allow
194  * the thread to continue running, which will lead to a crash if that thread then attempts to acquire the same semaphore.
195  * Without this, stress tests will crash.
196  */
197 void CancelSemaphoreOfThread(struct thread* thr) {
198     AssertSchedulerLockHeld();
199     assert(ThreadListContains(&thr->waiting_on_semaphore->waiting_list, thr));
200     ThreadListDelete(&thr->waiting_on_semaphore->waiting_list, thr);
201 }
202
203 int GetSemaphoreCount(struct semaphore* sem) {
204     return sem->current_count;
205 }

```

File: /fs/DS_Store

[binary]

File: /fs/filesystem.c

```

1
2 #include <common.h>
3 #include <vfs.h>
4 #include <errno.h>
5 #include <string.h>
6 #include <spinlock.h>
7 #include <irq.h>
8 #include <diskutil.h>
9 #include <semaphore.h>
10 #include <log.h>
11 #include <filesystem.h>
12 #include <heap.h>
13
14 #include <fs/demofs/demofs.h>
15
16 #define MAX_REGISTERED_FILESYSTEMS 8
17
18 struct filesystem {
19     char* name;
20     fs_mount_creator mount_creator;
21 };
22
23 static struct filesystem registered_filesystems[MAX_REGISTERED_FILESYSTEMS];
24 static int num_filesystems = 0;
25 static struct semaphore fs_table_lock;
26
27 void InitFilesystemTable(void) {
28     num_filesystems = 0;
29     fs_table_lock = CreateMutex("fs table");
30     RegisterFilesystem("demofs", DemofsMountCreator);
31 }
32
33 int RegisterFilesystem(char* fs_name, fs_mount_creator mount) {
34     if (fs_name == NULL || mount == NULL) {
35         return EINVAL;
36     }
37
38     struct filesystem fs;
39     fs.name = strdup_pageable(fs_name);
40     fs.mount_creator = mount;
41
42     int ret = 0;
43
44     AcquireMutex(fs_table_lock, -1);
45     if (num_filesystems < MAX_REGISTERED_FILESYSTEMS) {
46         registered_filesystems[num_filesystems++] = fs;
47     } else {
48         ret = EALREADY;
49     }
50     registered_filesystems[num_filesystems++] = fs;
51     ReleaseMutex(fs_table_lock);
52
53     if (ret != 0) {
54         FreeHeap(fs.name);
55     }
56
57     return ret;
58 }
59
60 int MountFilesystemForDisk(struct open_file* partition) {
61     AcquireMutex(fs_table_lock, -1);
62
63     struct open_file* fs = NULL;
64     for (int i = 0; i < num_filesystems; ++i) {
65         fs = NULL;
66         int res = registered_filesystems[i].mount_creator(partition, &fs);
67         if (res == 0) {
68             break;
69         }
70     }
71
72     ReleaseMutex(fs_table_lock);
73
74     if (fs == NULL) {
75         return ENODEV;
76     }
77
78     int res = VnodeOpCreate(partition->node, &fs->node, "fs", 0, 0);
79     if (res != 0) {
80         return res;
81     }
82
83     AddVfsMount(fs->node, GenerateNewMountedDiskName());
84     return 0;
85 }

```

File: /fs/demofs/demofs_private.h

```

#pragma once

#include <sys/types.h>
#include <common.h>
#include <vfs.h>
#include <transfer.h>

struct demofs {
    struct open_file* disk;
    ino_t root_inode;
};

#define MAX_NAME_LENGTH 24

#define INODE_TO_SECTOR(inode) (inode & 0xFFFFF)
#define INODE_IS_DIR(inode) (inode >> 31)
#define INODE_TO_DIR(inode) (inode | (1U << 31U))

int demofs_read_file(struct demofs* fs, ino_t file, uint32_t file_size, struct transfer* io);
int demofs_read_directory_entry(struct demofs* fs, ino_t directory, struct transfer* io);
int demofs_follow(struct demofs* fs, ino_t parent, ino_t* child, const char* name, uint32_t* file_length_out);

```

File: /fs/demofs/demofs_inodes.c

```

1
2 #include <common.h>
3 #include <errno.h>
4 #include <vfs.h>
5 #include <string.h>
6 #include <assert.h>
7 #include <panic.h>
8 #include <transfer.h>

```

```

9  #include <log.h>
10 #include <sys/types.h>
11 #include <dirent.h>
12 #include <fs/demofs/demofs_private.h>
13
14 #define SECTOR_SIZE 512
15
16 /*
17  * We are going to use the high bit of an inode ID to indicate whether or not
18  * we are talking about a directory or not (high bit set = directory).
19  *
20  * This allows us to, for example, easily catch ENOTDIR in demofs_follow.
21  *
22  * Remember to use INODE_TO_SECTOR. Note that inodes are only stored using 24 bits
23  * anyway.
24  */
25 int demofs_read_inode(struct demofs* fs, ino_t inode, uint8_t* buffer) {
26     struct transfer io = CreateKernelTransfer(buffer, SECTOR_SIZE, SECTOR_SIZE * INODE_TO_SECTOR(inode), TRANSFER_READ);
27     return ReadFile(fs->disk, &io);
28 }
29
30 int demofs_read_file(struct demofs* fs, ino_t file, uint32_t file_size_left, struct transfer* io) {
31     if (io->offset >= file_size_left) {
32         return 0;
33     }
34
35     file_size_left -= io->offset;
36
37     while (io->length_remaining != 0 && file_size_left != 0) {
38         int sector = file + io->offset / SECTOR_SIZE;
39         int sector_offset = io->offset % SECTOR_SIZE;
40
41         if (sector_offset == 0 && io->length_remaining >= SECTOR_SIZE && file_size_left >= SECTOR_SIZE) {
42             /*
43              * We have an aligned sector amount, so transfer it all directly,
44              * except for possible a few bytes at the end.
45              *
46              * ReadFile only allows lengths that are a multiple of the sector
47              * size, so round down to the nearest sector. The remainder must be
48              * kept track of so it can be added back on after the read.
49              */
50
51             int remainder = io->length_remaining % SECTOR_SIZE;
52             io->length_remaining -= remainder;
53
54             /*
55              * We need the disk offset, not the file offset.
56              * Ensure we move it back though afterwards.
57              */
58             int delta = sector * SECTOR_SIZE - io->offset;
59             io->offset += delta;
60
61             int status = ReadFile(fs->disk, io);
62             if (status != 0) {
63                 return status;
64             }
65
66             io->offset -= delta;
67             io->length_remaining = remainder;
68             file_size_left -= SECTOR_SIZE; // ???!
69
70         } else {
71             /*
72              * A partial sector transfer.
73              *
74              * We must read the sector into an internal buffer, and then copy a
75              * subsection of that to the return buffer.
76              */
77             uint8_t sector_buffer[SECTOR_SIZE];
78
79             /* Read the sector */
80             struct transfer temp_io = CreateKernelTransfer(sector_buffer, SECTOR_SIZE, sector * SECTOR_SIZE, TRANSFER_READ);
81
82             int status = ReadFile(fs->disk, &temp_io);
83             if (status != 0) {
84                 return status;
85             }
86
87             /* Transfer to the correct buffer */
88             size_t transfer_size = MIN(MIN(SECTOR_SIZE - (io->offset % SECTOR_SIZE), io->length_remaining), file_size_left);
89             PerformTransfer(sector_buffer + (io->offset % SECTOR_SIZE), io, transfer_size);
90             file_size_left -= transfer_size;
91         }
92     }
93
94     return 0;
95 }
96
97 int demofs_follow(struct demofs* fs, ino_t parent, ino_t* child, const char* name, uint32_t* file_length_out) {
98     assert(fs);
99     assert(fs->disk);
100     assert(child);
101     assert(name);
102     assert(file_length_out);
103     assert(SECTOR_SIZE % 32 == 0);
104
105     uint8_t buffer[SECTOR_SIZE];
106
107     if (strlen(name) > MAX_NAME_LENGTH) {
108         return ENAMETOOLONG;
109     }
110
111     if (!INODE_IS_DIR(parent)) {
112         return ENOTDIR;
113     }
114
115     /*
116      * The directory may contain many entries, so we need to iterate through them.
117      */
118     while (true) {
119         /*
120          * Grab the current entry.
121          */
122         int status = demofs_read_inode(fs, parent, buffer);
123         if (status != 0) {
124             return status;
125         }
126
127         /*
128          * Something went very wrong if the directory header is not present!
129          */
130         if (buffer[0] != 0xFF && buffer[0] != 0xFE) {
131             return EIO;
132         }
133
134         for (int i = 1; i < SECTOR_SIZE / 32; ++i) {
135             /*
136              * Check if there are no more names in the directory.
137              */
138             if (buffer[i * 32] == 0) {

```

```

139         return ENOENT;
140     }
141
142     /*
143     * Check if we've found the name.
144     */
145     if (!strcmp(name, (char*) buffer + i * 32, MAX_NAME_LENGTH)) {
146         /*
147         * If so, read the inode number and return it.
148         * Remember to add the directory flag if necessary.
149         */
150         ino_t inode = buffer[i * 32 + MAX_NAME_LENGTH + 4];
151         inode |= (ino_t) buffer[i * 32 + MAX_NAME_LENGTH + 5] << 8;
152         inode |= (ino_t) buffer[i * 32 + MAX_NAME_LENGTH + 6] << 16;
153
154         if (buffer[i * 32 + MAX_NAME_LENGTH + 7] & 1) {
155             /*
156             * This is a directory.
157             */
158             inode = INODE_TO_DIR(inode);
159             *file_length_out = 0;
160
161         } else {
162             /*
163             * This is a file.
164             */
165             uint32_t length = buffer[i * 32 + MAX_NAME_LENGTH];
166             length |= (uint32_t) buffer[i * 32 + MAX_NAME_LENGTH + 1] << 8;
167             length |= (uint32_t) buffer[i * 32 + MAX_NAME_LENGTH + 2] << 16;
168             length |= (uint32_t) buffer[i * 32 + MAX_NAME_LENGTH + 3] << 24;
169
170             *file_length_out = length;
171         }
172
173         *child = inode;
174         return 0;
175     }
176
177     /*
178     * Now we need to move on to the next entry if there is one.
179     */
180     if (buffer[0] == 0xFF) {
181         /* No more entries. */
182         return ENOENT;
183     } else if (buffer[0] == 0xFE) {
184         /*
185         * There is another entry, so read its inode and keep the loop going
186         */
187         parent = buffer[1];
188         parent |= (ino_t) buffer[2] << 8;
189         parent |= (ino_t) buffer[3] << 16;
190
191         /*
192         * Add the directory bit to the inode number as it should be a directory.
193         */
194         parent = INODE_TO_DIR(parent);
195     } else {
196         /*
197         * Something went very wrong if the directory header is not present!
198         */
199         return EIO;
200     }
201
202     }
203
204 }
205
206
207 int demofs_read_directory_entry(struct demofs fs, ino_t directory, struct transfer* io) {
208     if (!INODE_IS_DIR(directory)) {
209         return ENOTDIR;
210     }
211
212     assert(SECTOR_SIZE % 32 == 0);
213     uint8_t buffer[SECTOR_SIZE];
214
215     struct dirent dir;
216
217     if (io->offset % sizeof(struct dirent) != 0) {
218         return EINVAL;
219     }
220
221     int entry_number = io->offset / sizeof(struct dirent);
222
223     /*
224     * Each directory inode contains 31 files, and a pointer to the next directory entry.
225     * Add 1 to the offset to skip past the header.
226     */
227     int indirections = entry_number / 31;
228     int offset = entry_number % 31 + 1;
229
230     /*
231     * Get the correct inode
232     */
233     int status = 0;
234     ino_t current_inode = directory;
235     for (int i = 0; i < indirections; ++i) {
236         status = demofs_read_inode(fs, current_inode, buffer);
237         if (status != 0) {
238             return status;
239         }
240
241         /*
242         * Check for end of directory.
243         */
244         if (buffer[0] == 0xFF) {
245             return 0;
246         }
247
248         if (buffer[0] != 0xFE) {
249             return EIO;
250         }
251
252         /*
253         * Get the next in the chain.
254         */
255         current_inode = buffer[1];
256         current_inode |= (ino_t) buffer[2] << 8;
257         current_inode |= (ino_t) buffer[3] << 16;
258     }
259
260     status = demofs_read_inode(fs, current_inode, buffer);
261     if (status != 0) {
262         return status;
263     }
264
265     /*
266     * Check if we've gone past the end of the directory.
267     */
268     if (buffer[offset * 32] == 0) {

```

```

269     return 0;
270 }
271
272 /*
273  * strncpy is a bit iffy, so let's just play it safe.
274  */
275 char name[MAX_NAME_LENGTH + 2];
276 memset(name, 0, MAX_NAME_LENGTH + 2);
277 strncpy(name, (char*) buffer + offset * 32, MAX_NAME_LENGTH);
278 strcpy(dir.d_name, name);
279
280 dir.d_namlen = strlen(name);
281
282 ino_t inode = buffer[offset * 32 + MAX_NAME_LENGTH + 4];
283 inode |= (ino_t) buffer[offset * 32 + MAX_NAME_LENGTH + 5] << 8;
284 inode |= (ino_t) buffer[offset * 32 + MAX_NAME_LENGTH + 6] << 16;
285
286 dir.d_ino = inode;
287 dir.d_type = INODE_IS_DIR(inode) ? DT_DIR : DT_REG;
288
289 /* Perform the transfer to the correct location */
290 return PerformTransfer(&dir, io, sizeof(struct dirent));
291 }

```

File: ./fs/demofs/demofs.h

```
#pragma once
```

```
#include <vfs.h>
```

```
#include <common.h>
```

```
#include <filesystem.h>
```

```
int DemofsMountCreator(struct open_file* raw_device, struct open_file** out);
```

File: ./fs/demofs/demofs_vnode.c

```

1
2 #include <heap.h>
3 #include <log.h>
4 #include <assert.h>
5 #include <errno.h>
6 #include <fcntl.h>
7 #include <vfs.h>
8 #include <transfer.h>
9 #include <string.h>
10 #include <dirent.h>
11 #include <sys/stat.h>
12 #include <sys/types.h>
13 #include <fs/demofs/demofs_private.h>
14
15 struct vnode_data {
16     ino_t inode;
17     struct demofs fs;
18     uint32_t file_length;
19     bool directory;
20 };
21
22 static int CheckOpen(struct vnode*, const char* name, int flags) {
23     if (strlen(name) >= MAX_NAME_LENGTH) {
24         return ENAMETOOLONG;
25     }
26
27     if ((flags & O_ACCMODE) == O_WRONLY || (flags & O_ACCMODE) == O_RDWR) {
28         return EROFS;
29     }
30
31     return 0;
32 }
33
34 static int Ioctl(struct vnode*, int, void*) {
35     return EINVAL;
36 }
37
38 static int Read(struct vnode* node, struct transfer* io) {
39     struct vnode_data* data = node->data;
40     if (data->directory) {
41         return demofs_read_directory_entry(&data->fs, data->inode, io);
42     } else {
43         return demofs_read_file(&data->fs, data->inode, data->file_length, io);
44     }
45 }
46
47 static int Write(struct vnode*, struct transfer*) {
48     return EROFS;
49 }
50
51 static int Create(struct vnode*, struct vnode**, const char*, int, mode_t) {
52     return EROFS;
53 }
54
55 static int Truncate(struct vnode*, off_t) {
56     return EROFS;
57 }
58
59 static int Close(struct vnode* node) {
60     FreeHeap(node->data);
61     return 0;
62 }
63
64 static struct vnode* CreateDemoFsVnode(ino_t, off_t);
65
66 static int Follow(struct vnode* node, struct vnode** out, const char* name) {
67     struct vnode_data* data = node->data;
68     if (data->directory) {
69         ino_t child_inode;
70         uint32_t file_length;
71
72         int status = demofs_follow(&data->fs, data->inode, &child_inode, name, &file_length);
73         if (status != 0) {
74             return status;
75         }
76
77         /*
78          * TODO: return existing vnode if someone opens the same file twice...
79          */
80
81         struct vnode* child_node = CreateDemoFsVnode(child_inode, file_length);
82         struct vnode_data* child_data = AllocHeap(sizeof(struct vnode_data));
83         child_data->inode = child_inode;
84         child_data->fs = data->fs;
85         child_data->file_length = file_length;
86         child_data->directory = INODE_IS_DIR(child_inode);
87         child_node->data = child_data;

```

```

88
89     *out = child_node;
90
91     return 0;
92
93 } else {
94     return ENOTDIR;
95 }
96 }
97
98 static const struct vnode_operations dev_ops = {
99     .check_open   = CheckOpen,
100    .ioctl        = Ioctl,
101    .read         = Read,
102    .write        = Write,
103    .close        = Close,
104    .truncate     = Truncate,
105    .create       = Create,
106    .follow       = Follow,
107 };
108
109 static struct vnode* CreateDemoFsVnode(ino_t inode, off_t size) {
110     return CreateVnode(dev_ops, (struct stat) {
111         .st_mode = (INODE_IS_DIR(inode) ? S_IFDIR : S_IFREG) | S_IRWXU | S_IRWXG | S_IRWXO,
112         .st_nlink = 1,
113         .st_size = size,
114         .st_ino = inode,
115         .st_blksize = 512, // the 'efficient' size
116     });
117 }
118
119 static int CheckForDemofsSignature(struct open_file* raw_device) {
120     struct stat st = raw_device->node->stat;
121
122     uint8_t* buffer = AllocHeapEx(st.st_blksize, HEAP_ALLOW_PAGING);
123     struct transfer io = CreateKernelTransfer(buffer, st.st_blksize, 8 * st.st_blksize, TRANSFER_READ);
124     int res = ReadFile(raw_device, &io);
125     if (res != 0) {
126         FreeHeap(buffer);
127         return ENOTSUP;
128     }
129
130     /*
131     * Check for the DemoFS signature.
132     */
133     if (buffer[0] != 'D' || buffer[1] != 'E' || buffer[2] != 'M' || buffer[3] != 'O') {
134         FreeHeap(buffer);
135         return ENOTSUP;
136     }
137
138     return 0;
139 }
140
141 int DemofsMountCreator(struct open_file* raw_device, struct open_file** out) {
142     int sig_check = CheckForDemofsSignature(raw_device);
143     if (sig_check != 0) {
144         return sig_check;
145     }
146
147     struct vnode* node = CreateDemoFsVnode(9 | (1 << 31), 0);
148     struct vnode_data* data = AllocHeap(sizeof(struct vnode_data));
149
150     data->fs.disk = raw_device;
151     data->fs.root_inode = 9 | (1 << 31);
152     data->inode = 9 | (1 << 31); /* root directory inode */
153     data->file_length = 0; /* root directory has no length */
154     data->directory = true;
155
156     node->data = data;
157
158     *out = CreateOpenFile(node, 0, 0, true, false);
159     return 0;
160 }

```

File: ./thread/thread.c

```

1
2 #include <cpu.h>
3 #include <thread.h>
4 #include <spinlock.h>
5 #include <heap.h>
6 #include <assert.h>
7 #include <timer.h>
8 #include <string.h>
9 #include <irq.h>
10 #include <log.h>
11 #include <errno.h>
12 #include <virtual.h>
13 #include <panic.h>
14 #include <common.h>
15 #include <threadlist.h>
16 #include <avl.h>
17 #include <priorityqueue.h>
18 #include <progload.h>
19 #include <semaphore.h>
20 #include <process.h>
21
22 static struct thread_list ready_list;
23 static struct spinlock scheduler_lock;
24 static struct spinlock innermost_lock;
25
26 /*
27 * Local fixed sized arrays and variables need to fit on the kernel stack.
28 * Allocate at least 8KB (depending on the system page size).
29 *
30 * Please note that overflowing the kernel stack into non-paged memory will lead to
31 * an immediate and unrecoverable crash on most systems.
32 */
33 #define DEFAULT_KERNEL_STACK_KB BytesToPages(1024 * 16) * ARCH_PAGE_SIZE / 1024
34
35 /*
36 * The user stack is allocated as needed - this is the maximum size of the stack in
37 * user virtual memory. (However, a larger max stack means more page tables need to be
38 * allocated to store it - even if there are no actual stack pages in yet).
39 *
40 * On x86, allocating a 4MB region only requires one page table, hence we'll use that.
41 */
42 #define USER_STACK_MAX_SIZE BytesToPages(1024 * 1024 * 4) * ARCH_PAGE_SIZE
43
44 #define TIMESLICE_LENGTH_MS 50
45
46 /*
47 * Kernel stack overflow normally results in a total system crash/reboot because
48 * fault handlers will not work (they push data to a non-existent stack!).
49 *
50 * We will fill pages at the end of the stack with a certain value (CANARY_VALUE),

```



```

52 * and then we can check if they have been modified. If they are, we will throw a
53 * somewhat nicer error than a system reboot.
54 *
55 * Note that we can still overflow 'badly' if someone makes an allocation on the
56 * stack which is larger than the remaining space on the stack and the canary size
57 * combined.
58 *
59 * If the canary page is only partially used for the canary, the remainder of the
60 * page is able to be used normally.
61 */
62 #ifdef NDEBUG
63 #define NUM_CANARY_PAGES 0
64 #else
65 #define NUM_CANARY_BYTES (1024 * 8)
66 #define NUM_CANARY_PAGES BytesToPages(NUM_CANARY_BYTES)
67 #define CANARY_VALUE 0x8BADF00D
68
69 static void CheckCanary(size_t canary_base) {
70     uint32_t* canary_ptr = (uint32_t*) canary_base;
71
72     for (size_t i = 0; i < NUM_CANARY_BYTES / sizeof(uint32_t); ++i) {
73         if (*canary_ptr++ != CANARY_VALUE) {
74             Panic("PANIC_CANARY_DIED");
75         }
76     }
77 }
78
79 static void CreateCanary(size_t canary_base) {
80     uint32_t* canary_ptr = (uint32_t*) canary_base;
81
82     for (size_t i = 0; i < NUM_CANARY_BYTES / sizeof(uint32_t); ++i) {
83         *canary_ptr++ = CANARY_VALUE;
84     }
85 }
86
87 #endif
88
89
90 /*
91 * Allocates a new page-aligned stack for a kernel thread, and returns
92 * the address of either the top of the stack (if it grows downward),
93 * or the bottom (if it grows upward).
94 */
95 static void CreateKernelStacks(struct thread* thr, int kernel_stack_kb) {
96     int total_bytes = (BytesToPages(kernel_stack_kb * 1024) + NUM_CANARY_PAGES) * ARCH_PAGE_SIZE;
97
98     size_t stack_bottom = MapVirt(0, 0, total_bytes, VM_READ | VM_WRITE | VM_LOCK, NULL, 0);
99     size_t stack_top = stack_bottom + total_bytes;
100
101 #ifdef NDEBUG
102     thr->canary_position = stack_bottom;
103     CreateCanary(stack_bottom);
104 #endif
105     thr->kernel_stack_top = stack_top;
106     thr->kernel_stack_size = total_bytes;
107
108     thr->stack_pointer = ArchPrepareStack(thr->kernel_stack_top);
109 }
110
111 static size_t CreateUserStack(int size) {
112     /*
113     * All user stacks share the same area of virtual memory, but have different
114     * mappings to physical memory.
115     */
116
117     int total_bytes = BytesToPages(size) * ARCH_PAGE_SIZE;
118     size_t stack_base = ARCH_USER_STACK_LIMIT - total_bytes;
119     size_t actual_base = MapVirt(0, stack_base, total_bytes, VM_READ | VM_WRITE | VM_USER | VM_LOCAL, NULL, 0);
120
121     assert(stack_base == actual_base);
122     (void) actual_base; // The assert gets taken out on release mode, so make the compiler happy
123
124     return ARCH_USER_STACK_LIMIT;
125 }
126
127 /**
128 * Blocks the currently executing thread (no effect will happen until the IRQ goes below IRQ_SCHEDULER).
129 * The scheduler lock must be held.
130 */
131
132 void BlockThread(int reason) {
133     AssertSchedulerLockHeld();
134     assert(reason != THREAD_STATE_READY && reason != THREAD_STATE_RUNNING);
135     assert(GetCpu() != NULL && GetCpu()->current_thread != NULL);
136     assert(GetCpu()->current_thread->state == THREAD_STATE_RUNNING);
137     GetCpu()->current_thread->state = reason;
138     PostponeScheduleUntilStandardIrql();
139 }
140
141 void UnblockThread(struct thread* thr) {
142     AssertSchedulerLockHeld();
143     if (thr->state == THREAD_STATE_WAITING_FOR_SEMAPHORE_WITH_TIMEOUT) {
144         CancelSemaphoreOfThread(thr);
145     }
146     ThreadListInsert(&ready_list, thr);
147     if (thr->priority < GetThread()->priority) {
148         PostponeScheduleUntilStandardIrql();
149     }
150 }
151
152 void UpdateThreadTimeUsed(void) {
153     static uint64_t prev_time = 0;
154
155     uint64_t time = GetSystemTimer();
156     uint64_t time_elapsed = GetSystemTimer() - prev_time;
157     prev_time = time;
158
159     GetThread()->time_used += time_elapsed;
160 }
161
162 static int GetNextThreadId(void) {
163     static struct spinlock thread_id_lock;
164     static int next_thread_id = 0;
165     static bool initialised = false;
166
167     if (!initialised) {
168         initialised = true;
169         InitSpinlock(&thread_id_lock, "thread id", IRQL_SCHEDULER);
170     }
171
172     AcquireSpinlockIrql(&thread_id_lock);
173     int result = next_thread_id++;
174     ReleaseSpinlockIrql(&thread_id_lock);
175     return result;
176 }
177
178 struct thread* CreateThreadEx(void(*entry_point)(void*), void* argument, struct vas* vas, const char* name, struct process* prcss, int policy, int priority, int
179 struct thread* thr = AllocHeap(sizeof(struct thread));
180 thr->argument = argument;
181 thr->initial_address = entry_point;

```

```

182 thr->state = THREAD_STATE_READY;
183 thr->time_used = 0;
184 thr->name = strdup(name);
185 thr->priority = priority;
186 thr->needs_termination = false;
187 thr->time_used = false;
188 thr->waiting_on_semaphore = NULL;
189 thr->schedule_policy = policy;
190 thr->timeslice_expiry = GetSystemTimer() + TIMESLICE_LENGTH_MS;
191 thr->vas = vas;
192 thr->thread_id = GetNextThreadId();
193 CreateKernelStacks(thr, kernel_stack_kb == 0 ? DEFAULT_KERNEL_STACK_KB : 0);
194
195 if (prcss != NULL) {
196     AddThreadToProcess(prcss, thr);
197 } else {
198     thr->process = NULL;
199 }
200
201 LockScheduler();
202 ThreadListInsert(&ready_list, thr);
203 UnlockScheduler();
204
205 return thr;
206 }
207
208 struct thread* CreateThread(void(*entry_point)(void*), void* argument, struct vas* vas, const char* name) {
209     return CreateThreadEx(
210         entry_point, argument, vas, name, GetProcess(), SCHEDULE_POLICY_FIXED, FIXED_PRIORITY_KERNEL_NORMAL, 0
211     );
212 }
213
214 struct thread* GetThread(void) {
215     return GetCpu()->current_thread;
216 }
217
218 static void UpdateTimesliceExpiry(void) {
219     struct thread* thr = GetThread();
220     thr->timeslice_expiry = GetSystemTimer() + (thr->priority == 255 ? 0 : (20 + thr->priority / 4) * 1000000ULL);
221 }
222
223 void ThreadExecuteInUsermode(void* arg) {
224     struct thread* thr = GetThread();
225
226     int res = CopyProgramLoaderIntoAddressSpace();
227     if (res != 0) {
228         LogDeveloperWarning("COULDNT LOAD PROGRAM LOADER!\n");
229         TerminateThread(thr);
230     }
231
232     size_t user_stack = CreateUserStack(USER_STACK_MAX_SIZE);
233
234     LockScheduler();
235     thr->stack_pointer = user_stack;
236     UnlockScheduler();
237
238     ArchFlushTlb(GetVas());
239     ArchSwitchToUsermode(ARCH_PROG_LOADER_ENTRY, user_stack, arg);
240 }
241
242 struct process* CreateUsermodeProcess(struct process* parent, const char* filename) {
243     return CreateProcessWithEntryPoint(GetPid(parent), ThreadExecuteInUsermode, (void*) filename);
244 }
245
246 void ThreadInitialisationHandler(void) {
247     /*
248      * This normally happens in the schedule code, just after the call to ArchSwitchThread,
249      * but we forced ourselves to jump here instead, so we'd better do it now.
250      */
251     ReleaseSpinlockIrql(&innermost_lock);
252     UpdateTimesliceExpiry();
253
254     /*
255      * To get here, someone must have called thread_schedule(), and therefore
256      * the lock must have been held.
257      */
258     UnlockScheduler();
259
260     /* Anything else you might want to do should be done here... */
261
262     /* Go to the address the thread actually wanted. */
263     GetThread()->initial_address(GetThread()->argument);
264
265     /* The thread has returned, so just terminate it. */
266     TerminateThread(GetThread());
267
268     Panic(PANIC_IMPOSSIBLE_RETURN);
269 }
270
271 static int GetMinPriorityValueForPolicy(int policy) {
272     if (policy == SCHEDULE_POLICY_USER_HIGHER) return 50;
273     if (policy == SCHEDULE_POLICY_USER_NORMAL) return 100;
274     if (policy == SCHEDULE_POLICY_USER_LOWER) return 150;
275     return 0;
276 }
277
278 static int GetMaxPriorityValueForPolicy(int policy) {
279     if (policy == SCHEDULE_POLICY_FIXED) return 255;
280     return GetMinPriorityValueForPolicy(policy) + 100;
281 }
282
283 static void UpdatePriority(bool yielded) {
284     struct thread* thr = GetThread();
285     int policy = thr->schedule_policy;
286
287     if (policy != SCHEDULE_POLICY_FIXED) {
288         int new_val = thr->priority + (yielded ? -1 : 1);
289         if (new_val >= GetMinPriorityValueForPolicy(policy) && new_val <= GetMaxPriorityValueForPolicy(policy)) {
290             thr->priority = new_val;
291         }
292     }
293 }
294
295 void LockSchedulerX(void) {
296     AcquireSpinlockIrql(&scheduler_lock);
297 }
298
299 void UnlockSchedulerX(void) {
300     ReleaseSpinlockIrql(&scheduler_lock);
301 }
302
303 void AssertSchedulerLockHeld(void) {
304     assert(IsSpinlockHeld(&scheduler_lock));
305 }
306
307 __attribute__((returns_twice)) static void SwitchToNewTask(struct thread* old_thread, struct thread* new_thread) {
308     new_thread->state = THREAD_STATE_RUNNING;
309     ThreadListDeleteTop(&ready_list);
310 }

```

```

312
313 /*
314  * No IRQs allowed while this happens, as we need to protect the CPU structure.
315  * Only our CPU has access to it (as it is per-CPU), but if we IRQ and then someone
316  * calls GetCpu(), we'll be in a bit of strife.
317  */
318 struct cpu* cpu = GetCpu();
319 AcquireSpinlockIrql(&innermost_lock);
320
321 if (new_thread->vas != old_thread->vas) {
322     SetVas(new_thread->vas);
323 }
324
325 cpu->current_thread = new_thread;
326 cpu->current_vas = new_thread->vas;
327 ArchSwitchThread(old_thread, new_thread);
328
329 /*
330  * This code doesn't get called on the first time a thread gets run!! It jumps straight from
331  * ArchSwitchThread to ThreadInitialisationHandler!
332  */
333 ReleaseSpinlockIrql(&innermost_lock);
334
335 UpdateTimesliceExpiry();
336
337
338 static void ScheduleWithLockHeld(void) {
339     EXACT_IRQL(IRQL_SCHEDULER);
340     AssertSchedulerLockHeld();
341
342     struct thread* old_thread = GetThread();
343     struct thread* new_thread = ready_list.head;
344
345     if (old_thread == NULL) {
346         /*
347          * Multitasking not set up yet. Now check if someone has added a task that we can switch to.
348          * (If not, we keep waiting until they have, then we can start multitasking).
349          */
350         if (ready_list.head != NULL) {
351             /*
352              * We need a place where it can write the "old" stack pointer to.
353              */
354             struct thread dummy;
355             SwitchToNewTask(&dummy, new_thread);
356         }
357         return;
358     }
359
360     if (new_thread == old_thread || new_thread == NULL) {
361         /*
362          * Don't switch if there isn't anything else to switch to!
363          */
364         return;
365     }
366
367 #ifndef NDEBUG
368     CheckCanary(old_thread->canary_position);
369 #endif
370
371     bool yielded = old_thread->timeslice_expiry > GetSystemTimer();
372     UpdatePriority(yielded);
373     UpdateThreadTimeUsed();
374
375     /*
376      * Put the old task back on the ready list, but only if it didn't block / get suspended.
377      */
378     if (old_thread->state == THREAD_STATE_RUNNING) {
379         ThreadListInsert(&ready_list, old_thread);
380     }
381
382     SwitchToNewTask(old_thread, new_thread);
383 }
384
385 void Schedule(void) {
386     if (GetIrql() > IRQL_PAGE_FAULT) {
387         PostponeScheduleUntilStandardIrql();
388         return;
389     }
390
391     LockScheduler();
392     ScheduleWithLockHeld();
393     UnlockScheduler();
394
395     /**
396      * Used to allow TerminateThread() to kill a foreign process. This is because we can't just yank
397      * a thread off another list if it's blocked, as we don't know what list it's on. This way, we just
398      * signal that it needs terminating next time we allow it to run.
399      */
400     if (GetThread()->needs_termination) {
401         LogWriteSerial("Terminating a thread that was scheduled to die... stack at 0x%X\n", GetThread()->kernel_stack_top - GetThread()->kernel_stack_size);
402         TerminateThread(GetThread());
403         Panic(PANIC_IMPOSSIBLE_RETURN);
404     }
405 }
406
407 void InitScheduler() {
408     ThreadListInit(&ready_list, NEXT_INDEX_READY);
409     InitSpinlock(&scheduler_lock, "scheduler", IRQL_SCHEDULER);
410     InitSpinlock(&innermost_lock, "inner scheduler", IRQL_HIGH);
411 }
412
413 [[noreturn]] void StartMultitasking(void) {
414     InitIdle();
415     InitCleaner();
416
417     /*
418      * Once this is called, "the game is afoot!" and threads will start running.
419      */
420     Schedule();
421     Panic(PANIC_IMPOSSIBLE_RETURN);
422 }
423
424
425 /**
426  * Sets the priority and/or policy of a thread.
427  *
428  * @param policy The scheduling policy to use. Should be one of SCHEDULE_POLICY_FIXED, SCHEDULE_POLICY_USER_LOWER,
429  * SCHEDULE_POLICY_USER_NORMAL, SCHEDULE_POLICY_USER_HIGHER, or -1. Set to -1 to indicate that the
430  * policy should not be changed.
431  * @param priority The priority level to give the thread. Should be a value between 0 and 255, where 0 indicates the highest
432  * possible priority. A value of 255 indicates that the thread should only be run when the system is idle.
433  * If the thread's policy is SCHEDULE_POLICY_FIXED, then this value will get used directly. For other
434  * policies, the actual priority given to the thread may differ depending on the rules of the policy.
435  * Set to -1 to indicate that the policy should not be changed.
436  *
437  * @return Returns 0 on success, EINVAL if invalid arguments are given. If EINVAL is returned, no change will be made to the thread's
438  * policy or priority.
439  *
440  * @user This function may be used as a system call, as long as 'thr' points to a valid thread structure (which it should do,
441  * as the user will probably supply thread number, which the kernel then converts to address - or kernel may just make it

```

```

442 *         a 'current thread' syscall, in which case GetThread() will be valid.
443 *
444 * @maxirq1 IRQL_HIGH
445 */
446 int SetThreadPriority(struct thread* thr, int policy, int priority) {
447     if (priority < -1 || priority > 255) {
448         return EINVAL;
449     }
450     if (policy != -1 && policy != SCHEDULE_POLICY_FIXED && policy != SCHEDULE_POLICY_USER_LOWER && policy != SCHEDULE_POLICY_USER_NORMAL && policy != SCHEDULE_P
451         return EINVAL;
452     }
453     if (priority < GetMinPriorityValueForPolicy(policy)) {
454         priority = GetMinPriorityValueForPolicy(policy);
455     }
456     if (priority > GetMaxPriorityValueForPolicy(policy)) {
457         priority = GetMaxPriorityValueForPolicy(policy);
458     }
459     if (policy != -1) {
460         thr->scheduled_policy = policy;
461     }
462     if (priority != -1) {
463         thr->priority = priority;
464     }
465     return 0;
466 }
467
468 void AssignThreadToCpu(void) {
469     // NO-OP UNTIL SMP IMPLEMENTED
470 }
471
472 void UnassignThreadToCpu(void) {
473     // NO-OP UNTIL SMP IMPLEMENTED
474 }
475
476
477

```

File: ./thread/timer.c

```

1
2 #include <timer.h>
3 #include <spinlock.h>
4 #include <assert.h>
5 #include <cpu.h>
6 #include <panic.h>
7 #include <irq1.h>
8 #include <log.h>
9 #include <arch.h>
10 #include <thread.h>
11 #include <priorityqueue.h>
12 #include <threadlist.h>
13
14 static struct spinlock timer_lock;
15 static struct priority_queue sleep_queue;
16 static struct thread_list sleep_overflow_list;
17
18 #define SLEEP_QUEUE_LENGTH 32
19
20 static uint64_t system_time = 0;
21 static int sleep_wakeups_posted = 0;
22
23 void ReceivedTimer(uint64_t nanos) {
24     EXACT_IRQ1(IRQL_TIMER);
25
26     if (ArchGetCurrentCpuIndex() == 0) {
27         /*
28          * As we're in the timer handler, we know we already have IRQL_TIMER, and so we don't
29          * need to incur the additional overhead of raising and lowering.
30          */
31         AcquireSpinlockDirect(&timer_lock);
32         system_time += nanos;
33         ReleaseSpinlockDirect(&timer_lock);
34     }
35
36     /*
37      * Preempt the current thread if it has used up its timeslice.
38      */
39     struct thread* thr = GetThread();
40     if (thr != NULL && thr->timeslice_expiry != 0 && thr->timeslice_expiry <= system_time) {
41         PostponeScheduleUntilStandardIrql();
42     }
43
44     if (GetNumberInDeferQueue() < 8) {
45         DeferUntilIrql(IRQL_STANDARD, HandleSleepWakeups, (void*) &system_time);
46     }
47 }
48
49 uint64_t GetSystemTimer(void) {
50     MAX_IRQ1(IRQL_TIMER);
51
52     AcquireSpinlockIrql(&timer_lock);
53     uint64_t value = system_time;
54     ReleaseSpinlockIrql(&timer_lock);
55
56     return value;
57 }
58
59 void InitTimer(void) {
60     InitSpinlock(&timer_lock, "timer", IRQL_TIMER);
61     ThreadListInit(&sleep_overflow_list, NEXT_INDEX_SLEEP);
62     sleep_queue = PriorityQueueCreate(SLEEP_QUEUE_LENGTH, false, sizeof(struct thread*));
63 }
64
65 void QueueForSleep(struct thread* thr) {
66     AssertSchedulerLockHeld();
67
68     thr->timed_out = false;
69
70     if (PriorityQueueGetUsedSize(sleep_queue) == SLEEP_QUEUE_LENGTH) {
71         ThreadListInsert(&sleep_overflow_list, thr);
72     } else {
73         PriorityQueueInsert(sleep_queue, (void*) &thr, thr->sleep_expiry);
74     }
75 }
76
77 bool TryDequeueForSleep(struct thread* thr) {
78     AssertSchedulerLockHeld();
79
80     while (PriorityQueueGetUsedSize(sleep_queue) > 0) {
81         struct priority_queue_result res = PriorityQueuePeek(sleep_queue);
82         struct thread* top_thread = *((struct thread**) res.data);
83         PriorityQueuePop(sleep_queue);
84         if (top_thread == thr) {
85             return true;
86         }
87         ThreadListInsert(&sleep_overflow_list, top_thread);
88     }

```

```

89
90     struct thread* iter = sleep_overflow_list.head;
91     while (iter) {
92         if (iter == thr) {
93             ThreadListDelete(&sleep_overflow_list, iter);
94             return true;
95         } else {
96             iter = iter->next[NEXT_INDEX_SLEEP];
97         }
98     }
99 }
100
101 return false;
102 }
103
104 void HandleSleepWakeup(void* sys_time_ptr) {
105     MAX_IRQL(IRQL_PAGE_FAULT);
106
107     if (GetThread() == NULL) {
108         return;
109     }
110
111     LockScheduler();
112     if (sleep_wakeups_posted > 0) {
113         --sleep_wakeups_posted;
114     }
115
116     uint64_t system_time = *((uint64_t*) sys_time_ptr);
117
118     /*
119     * Wake up any sleeping tasks that need it.
120     */
121     while (PriorityQueueGetUsedSize(sleep_queue) > 0) {
122         struct priority_queue_result res = PriorityQueuePeek(sleep_queue);
123
124         /*
125         * Check if it needs waking.
126         */
127         if (res.priority <= system_time) {
128             struct thread* thr = *((struct thread**) res.data);
129             thr->timed_out = true;
130             PriorityQueuePop(sleep_queue);
131             UnblockThread(thr);
132         } else {
133             /*
134             * If this one doesn't need waking, none of the others will either.
135             */
136             break;
137         }
138     }
139 }
140
141 /*
142 * Check for any tasks that are asleep but on the overflow list. This is slow, but will
143 * only happen if we have more than 32 sleeping tasks, so it should normally not take any time.
144 */
145 struct thread* iter = sleep_overflow_list.head;
146 while (iter) {
147     if (iter->sleep_expiry <= system_time) {
148         ThreadListDelete(&sleep_overflow_list, iter);
149         iter->timed_out = true;
150         UnblockThread(iter);
151         iter = sleep_overflow_list.head;          // restart, as list changed
152     } else {
153         iter = iter->next[NEXT_INDEX_SLEEP];
154     }
155 }
156 }
157
158 UnlockScheduler();
159 }
160
161 /*
162 * Needs to be allowed at IRQL_PAGE_FAULT so the IDE driver can use it. Cannot be any higher, as otherwise
163 * the thread might not actually sleep (if IRQL_SCHEDULER or above, we won't actually do the 'blocked task switch'
164 * until it is released).
165 */
166 void SleepUntil(uint64_t system_time_ns) {
167     MAX_IRQL(IRQL_PAGE_FAULT);
168
169     if (system_time_ns < GetSystemTimer()) {
170         return;
171     }
172
173     LockScheduler();
174     GetThread()->sleep_expiry = system_time_ns;
175     QueueForSleep(GetThread());
176     BlockThread(THREAD_STATE_SLEEPING);
177     UnlockScheduler();
178 }
179
180 void SleepNano(uint64_t delta_ns) {
181     MAX_IRQL(IRQL_PAGE_FAULT);
182     SleepUntil(GetSystemTimer() + delta_ns);
183 }
184
185 void SleepMilli(uint32_t delta_ms) {
186     MAX_IRQL(IRQL_PAGE_FAULT);
187     SleepNano(((uint64_t) delta_ms) * 1000000ULL);
188 }

```

File: ./thread/process.c

```

1
2  /*
3   * thread/process.c - Processes
4   */
5
6  #include <arch.h>
7  #include <irq.h>
8  #include <thread.h>
9  #include <assert.h>
10 #include <virtual.h>
11 #include <sys/wait.h>
12 #include <sys/types.h>
13 #include <avl.h>
14 #include <panic.h>
15 #include <semaphore.h>
16 #include <filedes.h>
17 #include <spinlock.h>
18 #include <heap.h>
19 #include <process.h>
20 #include <log.h>
21 #include <linkedlist.h>
22
23 struct process {

```

```

24     pid_t pid;
25     struct vas* vas;
26     pid_t parent;
27     struct avl_tree* children;
28     struct avl_tree* threads;
29     struct semaphore* lock;
30     struct semaphore* killed_children_semaphore;
31     struct filedes_table* filedes_table;
32     int retv;
33     bool terminated;
34 };
35
36 struct process_table_node {
37     pid_t pid;
38     struct process* process;
39 };
40
41 static struct spinlock pid_lock;
42 static struct avl_tree* process_table;
43 static struct semaphore process_table_mutex;
44
45 static int ProcessTableComparator(void* a_, void* b_) {
46     struct process_table_node* a = a_;
47     struct process_table_node* b = b_;
48     return COMPARE_SIGN(a->pid, b->pid);
49 }
50
51 static pid_t AllocateNextPid(void) {
52     static pid_t next_pid = 1;
53
54     AcquireSpinlockIrq(&pid_lock);
55     pid_t pid = next_pid++;
56     ReleaseSpinlockIrq(&pid_lock);
57
58     return pid;
59 }
60
61 static int InsertIntoProcessTable(struct process* prcss) {
62     pid_t pid = AllocateNextPid();
63
64     AcquireMutex(&process_table_mutex, -1);
65
66     struct process_table_node* node = AllocHeap(sizeof(struct process_table_node));
67     node->pid = pid;
68     node->process = prcss;
69     AvlTreeInsert(process_table, (void*) node);
70
71     ReleaseMutex(&process_table_mutex);
72
73     return pid;
74 }
75
76 static void RemoveFromProcessTable(pid_t pid) {
77     AcquireMutex(&process_table_mutex, -1);
78
79     struct process_table_node dummy = {.pid = pid};
80     struct process_table_node* actual = AvlTreeGet(process_table, (void*) &dummy);
81     AvlTreeDelete(process_table, (void*) actual);
82     FreeHeap(actual); // this was allocated on 'InsertIntoProcessTable'
83
84     ReleaseMutex(&process_table_mutex);
85 }
86
87 void LockProcess(struct process* prcss) {
88     AcquireMutex(&prcss->lock, -1);
89 }
90
91 void UnlockProcess(struct process* prcss) {
92     ReleaseMutex(&prcss->lock);
93 }
94
95 void InitProcess(void) {
96     InitSpinlock(&pid_lock, "pid", IRQL_SCHEDULER);
97     process_table_mutex = CreateMutex("prcss table");
98     process_table = AvlTreeCreate();
99     AvlTreeSetComparator(process_table, ProcessTableComparator);
100 }
101
102 struct process* CreateProcess(pid_t parent_pid) {
103     EXACT_IRQL(IRQL_STANDARD);
104
105     struct process* prcss = AllocHeap(sizeof(struct process));
106
107     prcss->lock = CreateMutex("prcss");
108     prcss->vas = CreateVas();
109     prcss->parent = parent_pid;
110     prcss->children = AvlTreeCreate();
111     prcss->threads = AvlTreeCreate();
112     prcss->killed_children_semaphore = CreateSemaphore("killed children", SEM_BIG_NUMBER, SEM_BIG_NUMBER);
113     prcss->retv = 0;
114     prcss->terminated = false;
115     prcss->pid = InsertIntoProcessTable(prcss);
116     prcss->filedes_table = CreateFileDescriptorTable();
117
118     if (parent_pid != 0) {
119         struct process* parent = GetProcessFromPid(parent_pid);
120         LockProcess(parent);
121         AvlTreeInsert(parent->children, (void*) prcss);
122         UnlockProcess(parent);
123     }
124
125     return prcss;
126 }
127
128 void AddThreadToProcess(struct process* prcss, struct thread* thr) {
129     LockProcess(prcss);
130     AvlTreeInsert(prcss->threads, (void*) thr);
131     thr->process = prcss;
132     UnlockProcess(prcss);
133 }
134
135 struct process* ForkProcess(void) {
136     MAX_IRQL(IRQL_PAGE_FAULT);
137
138     LockProcess(GetProcess());
139
140     struct process* new_process = CreateProcess(GetProcess()->pid);
141     DestroyVas(new_process->vas);
142
143     // TODO: there are probably more things to copy over in the future (e.g. list of open file descriptors, etc.)
144     // the open files, etc.
145
146
147     // TODO: copy file descriptor table
148
149     new_process->vas = CopyVas();
150     //TODO: need to grab the first thread (I don't think we've ordered threads by thread id yet in the AVL)
151     // so will need to fix that first.
152     //CopyThreadToNewProcess(new_process, )
153     UnlockProcess(GetProcess());

```

```

154
155     return new_process;
156 }
157
158 /**
159  * Directly reaps a process.
160  */
161 static void ReapProcess(struct process* prcss) {
162     // TODO: there's more cleanup to be done here... ?
163     assert(prcss->vas != GetVas());
164
165     int res = DestroySemaphore(prcss->killed_children_semaphore, SEM_REQUIRE_FULL);
166     (void) res;
167     assert(res == 0);
168     DestroyVas(prcss->vas);
169     DestroyFileDescriptorTable(prcss->filedes_table);
170     RemoveFromProcessTable(prcss->pid);
171     if (prcss->parent != 0) {
172         struct process* parent = GetProcessFromPid(prcss->parent);
173         AvlTreeDelete(parent->children, prcss);
174     }
175     FreeHeap(prcss);
176 }
177
178 /**
179  * Recursively goes through the children of a process, reaping the first child that is able to be reaped.
180  * Depending on the value of 'target', it will either reap the first potential candidate, or a particular candidate.
181  *
182  * @param parent The process whose children we are looking through
183  * @param node The current subtree of the parent process' children tree
184  * @param target Either the process ID of the child to reap, or -1 to reap the first valid candidate.
185  * @param status If a child is reaped, its return value will be written here.
186  * @return The process ID of the reaped child, or 0 if no children are reaped.
187  */
188 static pid_t RecursivelyTryReap(struct process* parent, struct avl_node* node, pid_t target, int* status) {
189     if (node == NULL) {
190         return 0;
191     }
192
193     struct process* child = (struct process*) AvlTreeGetData(node);
194
195     LockProcess(child);
196
197     if (child->terminated && (child->pid == target || target == (pid_t) -1)) {
198         *status = child->retv;
199         pid_t pid = child->pid;
200         UnlockProcess(child); // needed in case someone is waiting on us, before our death
201         ReapProcess(child);
202         return pid;
203     }
204
205     UnlockProcess(child);
206
207     pid_t left_retv = RecursivelyTryReap(parent, AvlTreeGetLeft(node), target, status);
208     if (left_retv != 0) {
209         return left_retv;
210     }
211     return RecursivelyTryReap(parent, AvlTreeGetRight(node), target, status);
212 }
213
214
215 /**
216  * Changes the parent of a parentless process. Used to ensure the initial process can always reap orphaned
217  * processes.
218  */
219 static void AdoptOrphan(struct process* adopter, struct process* orphan) {
220     LockProcess(adopter);
221
222     orphan->parent = adopter->pid;
223     AvlTreeInsert(adopter->children, (void*) orphan);
224     ReleaseSemaphore(adopter->killed_children_semaphore);
225
226     UnlockProcess(adopter);
227 }
228
229 /**
230  * Recursively converts all child processes in a process' thread tree into zombie processes.
231  *
232  * @param node The subtree to start from. NULL is acceptable, and is the recursion base case.
233  */
234 static void RecursivelyMakeChildrenOrphans(struct avl_node* node) {
235     if (node == NULL) {
236         return;
237     }
238
239     RecursivelyMakeChildrenOrphans(AvlTreeGetLeft(node));
240     RecursivelyMakeChildrenOrphans(AvlTreeGetRight(node));
241     AdoptOrphan(GetProcessFromPid(1), AvlTreeGetData(node));
242 }
243
244 /**
245  * Recursively terminates all threads in a process' thread tree.
246  *
247  * @param node The subtree to start from. NULL is acceptable, and is the recursion base case.
248  */
249 static void RecursivelyKillRemainingThreads(struct avl_node* node) {
250     if (node == NULL) {
251         return;
252     }
253
254     RecursivelyKillRemainingThreads(AvlTreeGetLeft(node));
255     RecursivelyKillRemainingThreads(AvlTreeGetRight(node));
256
257     struct thread* victim = AvlTreeGetData(node);
258     if (victim->state != THREAD_STATE_TERMINATED && !victim->needs_termination) {
259         TerminateThread(victim);
260     }
261 }
262
263 /**
264  * Does all of the required operations to kill a process. This is run in its own thread, without an owning
265  * process, so that a process doesn't try to delete itself (and therefore delete its stack).
266  *
267  * @param arg The process to kill (needs to be cast to struct process*)
268  */
269 static void KillProcessHelper(void* arg) {
270     struct process* prcss = arg;
271
272     assert(GetProcess() == NULL);
273     assert(GetVas() != prcss->vas); // we should be on GetKernelVas()
274
275     RecursivelyKillRemainingThreads(AvlTreeGetRootNode(prcss->threads));
276     RecursivelyMakeChildrenOrphans(AvlTreeGetRootNode(prcss->children));
277
278     AvlTreeDestroy(prcss->threads);
279     AvlTreeDestroy(prcss->children);
280
281     DestroyVas(prcss->vas);
282
283     prcss->terminated = true;

```

```

284
285     if (prcss->parent == 0) {
286         ReapProcess(prcss);
287     } else {
288         struct process* parent = GetProcessFromPid(prcss->parent);
289         ReleaseSemaphore(parent->killed_children_semaphore);
290     }
291 }
292
293 TerminateThread(GetThread());
294
295
296 /**
297  * Deletes the current process and all its threads. Child processes have their parent switched to pid 1.
298  * If the process being deleted has a parent, then it becomes a zombie process until the parent reaps it
299  * If the process being deleted has no parent, it will be reaped and deallocated immediately.
300  *
301  * This function does not return.
302  *
303  * @param retv The return value the process being deleted will give.
304  */
305 void KillProcess(int retv) {
306     MAX_IRQL(IRQL_STANDARD);
307
308     struct process* prcss = GetProcess();
309     prcss->retv = retv;
310
311     /**
312      * Must run it in a different thread and process (a NULL process is fine), as it is going to kill all
313      * threads in the process, and the process itself.
314      */
315     CreateThreadEx(KillProcessHelper, (void*) prcss, GetKernelVas(), "process killer", NULL,
316                   SCHEDULE_POLICY_FIXED, FIXED_PRIORITY_KERNEL_HIGH, 0);
317
318     TerminateThread(GetThread());
319
320
321 /**
322  * Returns a pointer to the process that the thread currently running on this CPU belongs to. If there is no
323  * running thread (i.e. multitasking hasn't started yet), or the thread does not belong to a process, NULL
324  * is returned.
325  *
326  * @return The process of the current thread, if it exists, or NULL otherwise.
327  */
328 struct process* GetProcess(void) {
329     MAX_IRQL(IRQL_HIGH);
330     struct thread* thr = GetThread();
331     return thr == NULL ? NULL : thr->process;
332 }
333
334 struct process* CreateProcessWithEntryPoint(pid_t parent, void(*entry_point)(void*), void* args) {
335     EXACT_IRQL(IRQL_STANDARD);
336     struct process* prcss = CreateProcess(parent);
337     struct thread* thr = CreateThread(entry_point, args, prcss->vas, "prcssinit");
338     AddThreadToProcess(prcss, thr);
339     return prcss;
340 }
341
342 /**
343  * Returns the file descriptor table of the given process. Returns NULL if
344  * 'prcss' is null.
345  */
346 struct filedesc_table* GetFileDescriptorTable(struct process* prcss) {
347     if (prcss == NULL) {
348         return NULL;
349     }
350
351     return prcss->filedes_table;
352 }
353
354 /**
355  * Given a process id, returns the process object. Returns NULL for an invalid
356  * 'pid'.
357  */
358 struct process* GetProcessFromPid(pid_t pid) {
359     EXACT_IRQL(IRQL_STANDARD);
360
361     AcquireMutex(process_table_mutex, -1);
362
363     struct process_table_node dummy = {.pid = pid};
364     struct process_table_node* node = AvlTreeGet(process_table, (void*) &dummy);
365
366     ReleaseMutex(process_table_mutex);
367
368     return node == NULL ? NULL : node->process;
369 }
370
371 /**
372  * Returns the process id of a given process. If 'prcss' is null, 0 is returned.
373  */
374 pid_t GetPid(struct process* prcss) {
375     MAX_IRQL(IRQL_HIGH);
376     return prcss == NULL ? 0 : prcss->pid;
377 }
378
379 pid_t WaitProcess(pid_t pid, int* status, int flags) {
380     EXACT_IRQL(IRQL_STANDARD);
381
382     struct process* prcss = GetProcess();
383
384     pid_t result = 0;
385     int failed_reaps = 0;
386     while (result == 0) {
387         int res = AcquireSemaphore(prcss->killed_children_semaphore, (flags & WNOHANG) ? 0 : -1);
388         if (res != 0) {
389             break;
390         }
391         LockProcess(prcss);
392         result = RecursivelyTryReap(prcss, AvlTreeGetRootNode(prcss->children), pid, status);
393         UnlockProcess(prcss);
394         if (result == 0 && pid != (pid_t) -1) {
395             failed_reaps++;
396         }
397     }
398
399     /*
400      * Ensure that the next time we call WaitProcess(), we can immediately retry the reaps that
401      * we increased the semaphore for, but didn't actually reap on.
402      */
403     while (failed_reaps-- > 0) {
404         ReleaseSemaphore(prcss->killed_children_semaphore);
405     }
406
407     return result;
408 }

```


File: ./thread/progload.c

```
1
2 #include <thread.h>
3 #include <progload.h>
4 #include <assert.h>
5 #include <string.h>
6 #include <irq.h>
7 #include <fcntl.h>
8 #include <log.h>
9 #include <errno.h>
10 #include <virtual.h>
11 #include <panic.h>
12 #include <common.h>
13 #include <semaphore.h>
14 #include <sys/types.h>
15 #include <arch.h>
16 #include <vfs.h>
17
18 static size_t program_loader_addr;
19 static off_t program_loader_size;
20
21 void InitProgramLoader(void) {
22     struct open_file* file;
23     if (OpenFile("sys:/progload.exe", O_RDONLY, 0, &file)) {
24         PanicEx(PANIC_PROGRAM_LOADER, "program loader couldn't be loaded");
25     }
26
27     program_loader_size = file->node->stat.st_size;
28     program_loader_addr = MapVirt(0, 0, program_loader_size, VM_READ | VM_FILE, file, 0);
29 }
30
31 int CopyProgramLoaderIntoAddressSpace(void) {
32     size_t mem = MapVirt(0, ARCH_PROG_LOADER_BASE, program_loader_size, VM_READ | VM_EXEC | VM_WRITE | VM_USER | VM_LOCAL | VM_FIXED_VIRT, NULL, 0);
33     if (mem != ARCH_PROG_LOADER_BASE) {
34         return ENOMEM;
35     }
36
37     memcpy((void*) ARCH_PROG_LOADER_BASE, (void*) program_loader_addr, program_loader_size);
38     return 0;
39 }
```

File: ./thread/cleaner.c

```

1
2 /**
3  * thread/cleaner.c - Thread Termination Cleanup
4  *
5  * Threads are unable to delete their own stacks. Therefore, we have a separate thread which
6  * deletes the stacks (and any other leftover data) of threads that are marked as terminated.
7  */
8
9 #include <thread.h>
10 #include <virtual.h>
11 #include <threadlist.h>
12 #include <irq.h>
13 #include <heap.h>
14 #include <semaphore.h>
15 #include <log.h>
16 #include <panic.h>
17 #include <physical.h>
18
19 static struct thread_list terminated_list;
20 static struct semaphore* cleaner_semaphore;
21
22 static void DestroyThread(struct thread* thr) {
23     UnmapVirt(thr->kernel_stack_top - thr->kernel_stack_size, thr->kernel_stack_size);
24     FreeHeap(thr->name);
25     FreeHeap(thr);
26 }
27
28 static void CleanerThread(void*) {
29     while (true) {
30         AcquireSemaphore(cleaner_semaphore, -1);
31
32         LockScheduler();
33         struct thread* thr = terminated_list.head;
34         assert(thr != NULL);
35         ThreadListDeleteTop(&terminated_list);
36         UnlockScheduler();
37
38         DestroyThread(thr);
39     }
40 }
41
42 static void NotifyCleaner(void*) {
43     ReleaseSemaphore(cleaner_semaphore);
44 }
45
46 void TerminateThreadLockHeld(struct thread* thr) {
47     assert(thr != NULL);
48     EXACT_IRQL(IRQL_SCHEDULER);
49
50     if (thr == GetThread()) {
51         ThreadListInsert(&terminated_list, thr);
52
53         BlockThread(THREAD_STATE_TERMINATED);
54         DeferUntilIrql(IRQL_STANDARD, NotifyCleaner, NULL);
55     } else {
56         /*
57          * We can't terminate it directly, as it may be on any queue somewhere else. Instead, we
58          * will terminate it next time it is up for scheduling.
59          */
60         thr->needs_termination = true;
61     }
62 }
63
64 /**
65  * Terminates a thread. This function must not be called until after 'InitCleaner' has been called.
66  * The scheduler lock should not be already held.
67  *
68  * @param thr The thread to be terminated.
69  * @return This function does not return if thr == GetThread(), and returns void otherwise.
70  *
71  * @note MAX_IRQL(IRQL_SCHEDULER)
72  */
73 void TerminateThread(struct thread* thr) {
74     MAX_IRQL(IRQL_SCHEDULER);
75
76     LockScheduler();
77     TerminateThreadLockHeld(thr);
78     UnlockScheduler();
79
80     if (thr == GetThread()) {
81         Panic(PANIC_IMPOSSIBLE_RETURN);
82     }
83 }
84
85 /**
86  * Creates the cleaner thread. This must be called before any calls to 'TerminateThread' are made.
87  */
88 void InitCleaner(void) {
89     ThreadListInit(&terminated_list, NEXT_INDEX_TERMINATED);
90     cleaner_semaphore = CreateSemaphore("Cleaner", SEM_BIG_NUMBER, SEM_BIG_NUMBER);
91     CreateThread(CleanerThread, NULL, GetVas(), "cleaner");
92 }
93

```

File: ./thread/idle.c

```

1
2 /**
3  * thread/idle.c - System Idle Task
4  *
5  * A thread that is run if no other thread is available to run. The idle thread must therefore
6  * never block.
7  */
8
9 #include <arch.h>
10 #include <thread.h>
11 #include <virtual.h>
12 #include <irq.h>
13
14 static void IdleThread(void*) {
15     while (1) {
16         ArchStallProcessor();
17     }
18 }
19
20 void InitIdle(void) {
21     CreateThreadEx(IdleThread, NULL, GetVas(), "idle thread", NULL, SCHEDULE_POLICY_FIXED, FIXED_PRIORITY_IDLE, 0);
22 }
23

```

File: ./mem/swapfile.c

```

1 #include <virtual.h>
2 #include <vfs.h>
3 #include <fcntl.h>
4 #include <log.h>
5 #include <errno.h>
6 #include <panic.h>
7 #include <transfer.h>
8 #include <irq.h>
9 #include <virtual.h>
10 #include <physical.h>
11 #include <spinlock.h>
12 #include <arch.h>
13
14 static struct open_file* swapfile = NULL;
15 static struct spinlock swapfile_lock;
16 static uint8_t* swapfile_bitmap;
17 static int number_on_swapfile = 0;
18 static size_t num_swapfile_bitmap_entries = 0;
19
20 static int GetPagesRequiredForAllocationBitmap(void) {
21     uint64_t bits_per_page = ARCH_PAGE_SIZE * 8;
22     uint64_t accessible_per_page = ARCH_PAGE_SIZE * bits_per_page;
23     size_t max_swapfile_size = (GetTotalPhysKilobytes() * 1024) / 4 + (32 * 1024 * 1024);
24     return (max_swapfile_size * accessible_per_page - 1) / accessible_per_page;
25 }
26
27 static void SetupSwapfileBitmap(void) {
28     int num_pages_in_bitmap = GetPagesRequiredForAllocationBitmap();
29     num_swapfile_bitmap_entries = 8 * num_pages_in_bitmap * ARCH_PAGE_SIZE;
30     swapfile_bitmap = (uint8_t*) MapVirt(0, 0, num_pages_in_bitmap * ARCH_PAGE_SIZE, VM_READ | VM_WRITE | VM_LOCK, NULL, 0);
31 }
32
33 void InitSwapfile(void) {
34     if (OpenFile("swap:/", O_RDWR, 0, &swapfile)) {
35         Panic(PANIC_NO_FILESYSTEM);
36     }
37     InitSpinlock(&swapfile_lock, "swapfile", IRQL_SCHEDULER);
38     SetupSwapfileBitmap();
39 }
40
41
42 struct open_file* GetSwapfile(void) {
43     return swapfile;
44 }
45
46 static bool GetBitmapEntry(size_t index) {
47     return swapfile_bitmap[index / 8] & (1 << (index % 8));
48 }
49
50 static void SetBitmapEntry(size_t index, bool value) {
51     if (value) {
52         swapfile_bitmap[index / 8] |= 1 << (index % 8);
53     } else {
54         swapfile_bitmap[index / 8] &= ~(1 << (index % 8));
55     }
56 }
57
58 uint64_t AllocateSwapfileIndex(void) {
59     AcquireSpinlockIrql(&swapfile_lock);
60     ++number_on_swapfile;
61
62     for (size_t i = 0; i < num_swapfile_bitmap_entries; ++i) {
63         if (!GetBitmapEntry(i)) {
64             SetBitmapEntry(i, true);
65             ReleaseSpinlockIrql(&swapfile_lock);
66             return i;
67         }
68     }
69     Panic(PANIC_OUT_OF_SWAPFILE);
70 }
71
72
73 void DeallocateSwapfileIndex(uint64_t index) {
74     AcquireSpinlockIrql(&swapfile_lock);
75     --number_on_swapfile;
76     SetBitmapEntry(index, false);
77     ReleaseSpinlockIrql(&swapfile_lock);
78 }
79
80 int GetNumberOfPagesOnSwapfile(void) {
81     AcquireSpinlockIrql(&swapfile_lock);
82     int val = number_on_swapfile;
83     ReleaseSpinlockIrql(&swapfile_lock);
84     return val;
85 }

```

File: ./mem/heap.c

```

1 #include <common.h>
2 #include <assert.h>
3 #include <stdbool.h>
4 #include <panic.h>
5 #include <virtual.h>
6 #include <arch.h>
7 #include <string.h>
8 #include <spinlock.h>
9 #include <heap.h>
10 #include <log.h>
11 #include <semaphore.h>
12 #include <voidptr.h>
13 #include <irq.h>
14 #include <thread.h>
15
16 #define BOOTSTRAP_SIZE (1024 * 16)
17 #define MAX_EMERGENCY_BLOCKS 16
18
19 /*
20  * For non-pageable requests larger than this, we'll issue a warning to say that
21  * MapVirt is a much better choice.
22  */
23 #define WARNING_LARGE_REQUEST_SIZE (1024 * 3 + 512)
24
25 struct emergency_block {
26     uint8_t* address;
27     size_t size;
28     bool valid;
29 };
30
31 /**
32  * Used as a system block that can be used even before the virtual memory
33  * manager is setup.
34  */
35 static uint8_t bootstrap_memory_area[BOOTSTRAP_SIZE] __attribute__((aligned(ARCH_PAGE_SIZE)));
36
37 /*
38  * Used to give us memory when we are not allowed to fault (i.e. can't allocate

```

```

39  * virtual memory).
40  */
41  static struct emergency_block emergency_blocks[MAX_EMERGENCY_BLOCKS] = {
42      { .address = bootstrap_memory_area, .size = BOOTSTRAP_SIZE, .valid = true }
43  };
44
45  static struct semaphore* heap_lock;
46  static struct spinlock heap_spinlock;
47  static struct spinlock heap_locker_lock;
48  static struct thread* lock_entry_threads[2];
49
50  static bool LockHeap(bool paging) {
51      bool acquired = false;
52      struct thread* thr = GetThread();
53
54      AcquireSpinlockIrql(&heap_locker_lock);
55      if (lock_entry_threads[paging ? 1 : 0] != thr || lock_entry_threads[paging ? 1 : 0] == NULL) {
56          ReleaseSpinlockIrql(&heap_locker_lock);
57
58          if (paging) {
59              AcquireMutex(heap_lock, -1);
60          } else {
61              AcquireSpinlockIrql(&heap_spinlock);
62          }
63
64          lock_entry_threads[paging ? 1 : 0] = thr;
65          acquired = true;
66
67      } else {
68          ReleaseSpinlockIrql(&heap_locker_lock);
69      }
70
71      return acquired;
72  }
73
74  static void UnlockHeap(bool paging) {
75      lock_entry_threads[paging ? 1 : 0] = NULL;
76
77      if (paging) {
78          ReleaseMutex(heap_lock);
79      } else {
80          ReleaseSpinlockIrql(&heap_spinlock);
81      }
82  }
83
84  static void* AllocateFromEmergencyBlocks(size_t size) {
85      int smallest_block = -1;
86
87      for (int i = 0; i < MAX_EMERGENCY_BLOCKS; ++i) {
88          if (emergency_blocks[i].valid && emergency_blocks[i].size >= size) {
89              if (smallest_block == -1 || emergency_blocks[i].size < emergency_blocks[smallest_block].size) {
90                  smallest_block = i;
91              }
92          }
93      }
94
95      if (smallest_block == -1) {
96          Panic(PANIC_CANNOT_MALLOC_WITHOUT_FAULTING);
97      }
98
99      void* address = emergency_blocks[smallest_block].address;
100
101      emergency_blocks[smallest_block].address += size;
102      emergency_blocks[smallest_block].size -= size;
103
104      if (emergency_blocks[smallest_block].size < ARCH_PAGE_SIZE) {
105          emergency_blocks[smallest_block].valid = false;
106      }
107
108      return address;
109  }
110
111  /**
112   * This function needs to be called with the heap lock held.
113   */
114  static void AddBlockToBackupHeap(size_t size) {
115      UnlockHeap(false);
116      void* address = (void*) MapVirt(0, 0, size, VM_READ | VM_WRITE | VM_LOCK, NULL, 0);
117      LockHeap(false);
118
119      int index_of_smallest_block = 0;
120      for (int i = 0; i < MAX_EMERGENCY_BLOCKS; ++i) {
121          if (emergency_blocks[i].valid) {
122              if (emergency_blocks[i].size < emergency_blocks[index_of_smallest_block].size) {
123                  index_of_smallest_block = i;
124              }
125          } else {
126              emergency_blocks[i].valid = true;
127              emergency_blocks[i].size = size;
128              emergency_blocks[i].address = address;
129              return;
130          }
131      }
132
133      LogDeveloperWarning("losing 0x%X bytes due to strangeness with backup heap.\n", emergency_blocks[index_of_smallest_block].size);
134      LogDeveloperWarning("TODO: could probably just add this as a regular block to the regular heap.\n");
135
136      emergency_blocks[index_of_smallest_block].size -= size;
137      emergency_blocks[index_of_smallest_block].address = address;
138
139  }
140
141  static void RestoreEmergencyPages(void* context) {
142      (void) context;
143
144      EXACT_IRQL(IRQL_STANDARD);
145
146      /*
147       * TODO: maybe make this greedier (i.e. grab larger blocks), but also have a way for the PMM to ask for larger
148       * blocks back if needed. would still need to retain enough stashed away for PMM/VMM to use the heap, and would
149       * need to unlock the pages, and ensure that allocations from emergency blocks are done via smallest-fit (so large
150       * blocks aren't wasted).
151       */
152
153      size_t total_size = 0;
154      size_t largest_block = 0;
155
156      LockHeap(false);
157      for (int i = 0; i < MAX_EMERGENCY_BLOCKS; ++i) {
158          if (emergency_blocks[i].valid) {
159              size_t size = emergency_blocks[i].size;
160              total_size += size;
161              if (size > largest_block) {
162                  largest_block = size;
163              }
164          }
165      }
166
167      while (largest_block < BOOTSTRAP_SIZE / 2 || total_size < BOOTSTRAP_SIZE) {
168          AddBlockToBackupHeap(BOOTSTRAP_SIZE);
169      }

```

```

169     total_size += BOOTSTRAP_SIZE;
170     largest_block = BOOTSTRAP_SIZE;
171 }
172
173 UnlockHeap(false);
174
175
176 /**
177  * Represents a section of memory that is either allocated or free. The memory
178  * address it represents is itself, excluding the metadata at the start or end.
179  */
180 struct block {
181     /*
182      * The entire size of the block. The low 2 bits do not form part of the
183      * size, the low bit is set for allocated blocks, and the second lowest bit
184      * is indicates if it is on the swappable heap or not.
185      */
186     size_t size;
187
188     /*
189      * Only here on free blocks. Allocated blocks use this as the start of
190      * allocated memory.
191      */
192     struct block* next;
193     struct block* prev;
194
195     /*
196      * At position size - sizeof(size_t), there is the trailing size tag. There
197      * are no flags in the low bit of this value, unlike the heading size tag.
198      */
199 };
200
201 #ifndef NDEBUG
202 static int outstanding_allocations = 0;
203
204 int DbgGetOutstandingHeapAllocations(void) {
205     return outstanding_allocations;
206 }
207 #endif
208
209 /**
210  * Must be a power of 2.
211  */
212 #define ALIGNMENT 8
213
214 /**
215  * The amount of metadata at the start and end of allocated blocks. The next and free
216  * pointers in free blocks do not count.
217  */
218 #define METADATA_LEADING (sizeof(size_t))
219 #define METADATA_TRIALING (sizeof(size_t))
220 #define METADATA_TOTAL (METADATA_LEADING + METADATA_TRIALING)
221
222 #define MINIMUM_REQUEST_SIZE_INTERNAL (2 * sizeof(size_t))
223 #define TOTAL_NUM_FREE_LISTS 35
224
225 /**
226  * An array which holds the minimum allocation sizes that each free list can hold.
227  */
228 static size_t free_list_block_sizes[TOTAL_NUM_FREE_LISTS] = {
229     8,          16,          24,          32,
230     40,         48,          56,          64,
231     72,         80,          88,          96,
232     104,        112,         120,         128,
233     160,        192,         224,         256,
234     320,        384,         448,         512,
235     768,        1024,        1536,        2048,    // 28 [27]
236     1024 * 4,   1024 * 8,   1024 * 16,   1024 * 32,    // 32 [31]
237     1024 * 64,  1024 * 128, 1024 * 256,
238 };
239
240 /**
241  * Used to work out which free list a block should be in, when we are *reading*
242  * a block. This rounds the size *up*, meaning it cannot be used to insert a
243  * block into a list. NOT USED TO INSERT BLOCKS!!
244  */
245 static int GetSmallestListIndexThatFits(size_t size_without_metadata) {
246     int i = 0;
247     while (true) {
248         if (size_without_metadata <= free_list_block_sizes[i]) {
249             return i;
250         }
251         ++i;
252     }
253 }
254
255 /**
256  * Calculates which free list a block should be inserted in. This rounds down,
257  * and so it should not normally be used to look up where a block should be.
258  */
259 static int GetInsertionIndex(size_t size_without_metadata) {
260     assert(size_without_metadata >= free_list_block_sizes[0]);
261
262     /*
263      * We can't round down to the next one when it's the smallest possible size,
264      * so handle this case specially.
265      */
266     if (size_without_metadata == free_list_block_sizes[0]) {
267         return 0;
268     }
269
270     for (int i = 0; i < TOTAL_NUM_FREE_LISTS - 1; ++i) {
271         if (size_without_metadata <= free_list_block_sizes[i]) {
272             return i - 1;
273         }
274     }
275
276     Panic(PANIC_HEAP_REQUEST_TOO_LARGE);
277 }
278
279 /**
280  * Rounds up a user-supplied allocation size to the alignment. If it's less than
281  * the minimum size internally supported, we will round it up to that size.
282  */
283 static size_t RoundUpSize(size_t size) {
284     assert(size != 0);
285
286     if (size < MINIMUM_REQUEST_SIZE_INTERNAL) {
287         size = MINIMUM_REQUEST_SIZE_INTERNAL;
288     }
289
290     return (size + ALIGNMENT - 1) & ~(ALIGNMENT - 1);
291 }
292
293 static void MarkFree(struct block* block) {
294     block->size &= ~1;
295 }
296
297 static void MarkAllocated(struct block* block) {
298     block->size |= 1;

```

```

299 |
300 |
301 static bool IsAllocated(struct block* block) {
302     return block->size & 1;
303 |
304 |
305 static void MarkSwappability(struct block* block, int can_swap) {
306     if (can_swap)
307         block->size |= 2;
308     } else
309         block->size &= ~2;
310 |
311 |
312 |
313 static bool IsOnSwappableHeap(struct block* block) {
314     return block->size & 2;
315 |
316 |
317 /**
318  * Global arrays always initialise to zero (and therefore, to NULL).
319  * Entries in free lists must have a user allocated size GREATER OR EQUAL TO the
320  * size in free_list_block_sizes.
321  */
322 static struct block* _head_block[TOTAL_NUM_FREE_LISTS];
323 static struct block* _head_block_swappable[TOTAL_NUM_FREE_LISTS];
324 |
325 static struct block** GetHeap(bool swappable) {
326     return swappable ? _head_block_swappable : _head_block;
327 |
328 |
329 static struct block** GetHeapForBlock(struct block* block) {
330     return GetHeap(IsOnSwappableHeap(block));
331 |
332 |
333 |
334 /**
335  * Given a block, returns its total size, including metadata. This takes into
336  * account the flags on the size field and removes them from the return value.
337  */
338 static size_t GetSize(struct block* block) {
339     size_t size = block->size & ~3;
340     assert((((size_t*) block) + (size / sizeof(size_t)) - 1) == size);
341     return size;
342 |
343 |
344 /**
345  * Sets the *total* size of a given block. This does not do any checking, so the
346  * caller must be careful as setting the wrong size can corrupt the heap.
347  * Carries the swappability and allocation tags with it from the front to the
348  * back tags.
349  */
350 static void SetSizeTags(struct block* block, size_t size) {
351     block->size = (block->size & 3) | size;
352     (((size_t*) block) + (size / sizeof(size_t)) - 1) = size;
353 |
354 |
355 static void* GetSystemMemory(size_t size, int flags) {
356     if (flags & HEAP_NO_FAULT) {
357         if (flags & HEAP_ALLOW_PAGING) {
358             Panic(PANIC_CONFLICTING_ALLOCATION_REQUIREMENTS);
359         }
360         if (IsVirtInitialised()) {
361             DeferUntilIrql(IRQL_STANDARD, RestoreEmergencyPages, NULL);
362         }
363         return AllocateFromEmergencyBlocks(size);
364     }
365 |
366     return (void*) MapVirt(0, 0, size, VM_READ | VM_WRITE | (flags & HEAP_ALLOW_PAGING ? 0 : VM_LOCK), NULL, 0);
367 |
368 |
369 /**
370  * Allocates a new block from the system that is able to hold the amount of data
371  * specified. Also allocated enough memory for fenceposts on either side of the
372  * data, and sets up these fenceposts correctly.
373  */
374 static struct block* RequestBlock(size_t total_size, int flags) {
375     /*
376      * We need to add the extra bytes for fenceposts to be added. We must do
377      * this before we round up to the nearest areana size (if we did it after,
378      * it wouldn't be aligned anymore).
379      */
380     total_size += MINIMUM_REQUEST_SIZE_INTERNAL * 2;
381     total_size = (total_size + ARCH_PAGE_SIZE - 1) & ~(ARCH_PAGE_SIZE - 1);
382 |
383     if (!IsVirtInitialised()) {
384         flags |= HEAP_NO_FAULT;
385     }
386 |
387     struct block* block = (struct block*) GetSystemMemory(total_size, flags);
388     if (block == NULL) {
389         Panic(PANIC_OUT_OF_HEAP);
390     }
391 |
392     /*
393      * Set the metadata for both the fenceposts and the main data block.
394      * Keep in mind that total_size now includes the fencepost metadata (see top
395      * of function), so this sometimes needs to be subtracted off.
396      */
397     struct block* left_fence = block;
398     struct block* actual_block = (struct block*) (((size_t*) block) + MINIMUM_REQUEST_SIZE_INTERNAL / sizeof(size_t));
399     struct block* right_fence = (struct block*) (((size_t*) block) + (total_size - MINIMUM_REQUEST_SIZE_INTERNAL) / sizeof(size_t));
400 |
401     SetSizeTags(left_fence, MINIMUM_REQUEST_SIZE_INTERNAL);
402     SetSizeTags(actual_block, total_size - 2 * MINIMUM_REQUEST_SIZE_INTERNAL);
403     SetSizeTags(right_fence, MINIMUM_REQUEST_SIZE_INTERNAL);
404 |
405     actual_block->prev = NULL;
406     actual_block->next = NULL;
407 |
408     MarkAllocated(left_fence);
409     MarkAllocated(right_fence);
410     MarkFree(actual_block);
411 |
412     MarkSwappability(left_fence, flags & HEAP_ALLOW_PAGING);
413     MarkSwappability(right_fence, flags & HEAP_ALLOW_PAGING);
414     MarkSwappability(actual_block, flags & HEAP_ALLOW_PAGING);
415 |
416     return actual_block;
417 |
418 |
419 /**
420  * Removes a block from a free list. It needs to take in the exact free list's
421  * index (as opposed to calculating it itself), as this may be used halfway
422  * through allocations or deallocations where the block isn't yet in its correct
423  * block.
424  */
425 static void RemoveBlock(int free_list_index, struct block* block) {
426     struct block** head_list = GetHeapForBlock(block);
427 |
428     if (block->prev == NULL && block->next == NULL) {

```

```

429     assert(head_list[free_list_index] == block);
430     head_list[free_list_index] = NULL;
431
432     } else if (block->prev == NULL) {
433         head_list[free_list_index] = block->next;
434         block->next->prev = NULL;
435     }
436     } else if (block->next == NULL) {
437         block->prev->next = NULL;
438     }
439     } else {
440         block->prev->next = block->next;
441         block->next->prev = block->prev;
442     }
443 }
444
445 /**
446  * Adds a block to its appropriate free list. It also coalesces the block with
447  * surrounding free blocks if possible.
448  */
449 static struct block* AddBlock(struct block* block) {
450     size_t size = GetSize(block);
451     struct block* head_list = GetHeapForBlock(block);
452
453     int free_list_index = GetInsertionIndex(size - METADATA_TOTAL);
454
455     size_t prev_block_size = (((size_t*) block) - 1);
456     struct block* prev_block = (struct block*) (((size_t*) block) - prev_block_size / sizeof(size_t));
457     struct block* next_block = (struct block*) (((size_t*) block) + size / sizeof(size_t));
458
459     if (IsAllocated(prev_block) && IsAllocated(next_block)) {
460         /*
461          * Cannot coalesce here, so just add to the free list.
462          */
463         block->prev = NULL;
464         block->next = head_list[free_list_index];
465         if (block->next != NULL)
466             block->next->prev = block;
467     }
468     head_list[free_list_index] = block;
469     MarkFree(block);
470     return block;
471
472     } else if (IsAllocated(prev_block) && !IsAllocated(next_block)) {
473         /*
474          * Need to coalesce with the one on the right.
475          */
476         bool swappable = IsOnSwappableHeap(block);
477         assert(swappable == IsOnSwappableHeap(next_block));
478
479         RemoveBlock(GetInsertionIndex(GetSize(next_block) - METADATA_TOTAL), next_block);
480         SetSizeTags(block, size + GetSize(next_block));
481         block->prev = NULL;
482         block->next = NULL;
483         MarkFree(block);
484         MarkSwappability(block, swappable);
485
486         return AddBlock(block);
487     } else if (!IsAllocated(prev_block) && IsAllocated(next_block)) {
488         /*
489          * Need to coalesce with the one on the left.
490          */
491         bool swappable = IsOnSwappableHeap(block);
492         assert(swappable == IsOnSwappableHeap(prev_block));
493
494         RemoveBlock(GetInsertionIndex(GetSize(prev_block) - METADATA_TOTAL), prev_block);
495         SetSizeTags(prev_block, size + GetSize(prev_block));
496         prev_block->prev = NULL;
497         prev_block->next = NULL;
498         MarkFree(prev_block);
499         MarkSwappability(prev_block, swappable);
500
501         return AddBlock(prev_block);
502     } else {
503         /*
504          * Coalesce with blocks on both sides.
505          */
506         bool swappable = IsOnSwappableHeap(block);
507         assert(swappable == IsOnSwappableHeap(prev_block));
508         assert(swappable == IsOnSwappableHeap(next_block));
509
510         RemoveBlock(GetInsertionIndex(GetSize(prev_block) - METADATA_TOTAL), prev_block);
511         RemoveBlock(GetInsertionIndex(GetSize(next_block) - METADATA_TOTAL), next_block);
512
513         SetSizeTags(prev_block, size + GetSize(prev_block) + GetSize(next_block));
514         prev_block->prev = NULL;
515         prev_block->next = NULL;
516         MarkFree(prev_block);
517         MarkSwappability(prev_block, swappable);
518
519         return AddBlock(prev_block);
520     }
521 }
522
523 /**
524  * Allocates a block. The block to be allocated will be the first block in the
525  * given free list, and that free list must be non-empty, and be able to fit the
526  * requested size.
527  */
528 static struct block* AllocateBlock(struct block* block, int free_list_index, size_t user_requested_size) {
529     assert(block != NULL);
530
531     size_t total_size = user_requested_size + METADATA_TOTAL;
532     size_t block_size = GetSize(block);
533
534     assert(block_size >= total_size);
535
536     if (block_size - total_size < MINIMUM_REQUEST_SIZE_INTERNAL + METADATA_TOTAL) {
537         /*
538          * We can just remove from the list altogether if the sizes match up
539          * exactly, or if there would be so little left over that we can't form
540          * a new block.
541          */
542         RemoveBlock(free_list_index, block);
543
544         /*
545          * Prevent memory leak (from having a hole in memory), but do it after
546          * removing the block, as this may change the list it needs to be in,
547          * and RemoveBlock will not like that.
548          */
549         SetSizeTags(block, block_size);
550         MarkAllocated(block);
551         return block;
552     } else {
553         /*
554          * We must split the block into two. If no list change is needed, we can
555          * leave the 'leftover' parts in the list as is (just fixing up the size
556          * tags), and then return the new block.
557          */
558     }

```

```

559     */
560
561     RemoveBlock(free_list_index, block);
562
563     size_t leftover = block_size - total_size;
564     SetSizeTags(block, leftover);
565
566     struct block* allocated_block = (struct block*) (((size_t*) block) + (leftover / sizeof(size_t)));
567     SetSizeTags(allocated_block, total_size);
568
569     /*
570     * We are giving it new tags, so must set this correctly.
571     */
572     MarkSwappability(allocated_block, IsOnSwappableHeap(block));
573
574     /*
575     * Must be done before we try to move around the leftovers (or else it
576     * will coalesce back into one block).
577     */
578     MarkAllocated(allocated_block);
579
580     /*
581     * We need to remove the leftover block from this list, and add it to the
582     * correct list.
583     */
584     MarkFree(block);
585     AddBlock(block);
586     return allocated_block;
587 }
588 }
589
590 /**
591 * Allocates a block that can fit the user requested size. It will request new
592 * memory from the system if required. If it returns NULL, then there is not
593 * enough memory of the system to satisfy the request.
594 */
595 static struct block* FindBlock(size_t user_requested_size, int flags) {
596     struct block* head_list = GetHeap(flags & HEAP_ALLOW_PAGING);
597
598     int min_index = GetSmallestListIndexThatFits(user_requested_size);
599     for (int i = min_index; i < TOTAL_NUM_FREE_LISTS; ++i) {
600         if (head_list[i] != NULL) {
601             return AllocateBlock(head_list[i], i, user_requested_size);
602         }
603     }
604
605     /*
606     * If we can't find a block that will fit, then we must allocate memory.
607     * Round up to the next block size, so we ensure we are in the next bucket.
608     * This avoids an issue if e.g. a user requests 2.1KB, and we allocate 3.9KB
609     * and it goes in the wrong bucket due to the two different indexes used.
610     */
611     size_t total_size = free_list_block_sizes[min_index + 1] + METADATA_TOTAL;
612     struct block* sys_block = RequestBlock(total_size, flags);
613
614     /*
615     * Put the new memory in the free list (which ought to be empty, as wouldn't
616     * need to request new memory otherwise). Then we can allocate the block.
617     */
618     int sys_index = GetInsertionIndex(GetSize(sys_block) - METADATA_TOTAL);
619     assert(head_list[sys_index] == NULL);
620     head_list[sys_index] = sys_block;
621     return AllocateBlock(head_list[sys_index], sys_index, user_requested_size);
622 }
623
624 /**
625 * Allocates memory on the heap. Unless you *really* know what you're doing, you
626 * should always pass HEAP_NO_FAULT. AllocHeap passes this automatically, but
627 * this one doesn't (in case you want to allocate from the pagable pool).
628 */
629 void* AllocHeapEx(size_t size, int flags) {
630     MAX_IRQL(IRQL_SCHEDULER);
631     if (flags & HEAP_ALLOW_PAGING) {
632         EXACT_IRQL(IRQL_STANDARD);
633     }
634
635     if (size == 0) {
636         return NULL;
637     }
638     if (size >= WARNING_LARGE_REQUEST_SIZE && ((flags & HEAP_ALLOW_PAGING) == 0 || (flags & HEAP_NO_FAULT) != 0)) {
639         LogDeveloperWarning("AllocHeapEx called with allocation of size 0x%X. You should consider using MapVirt.\n", size);
640     }
641     if (flags == 0) {
642         LogDeveloperWarning("AllocHeapEx called with flags = 0. You probably meant to pass either HEAP_ALLOW_PAGING,"
643             "or HEAP_NO_FAULT. Passing neither is valid and it puts it on the locked heap, but allocation"
644             "may cause faults. This is unlikely to be what you want.");
645     }
646
647     bool acquired = LockHeap(flags & HEAP_ALLOW_PAGING);
648     size = RoundUpSize(size);
649     struct block* block = FindBlock(size, flags);
650
651 #ifndef NDEBUG
652     outstanding_allocations++;
653 #endif
654
655     if (acquired) {
656         UnlockHeap(flags & HEAP_ALLOW_PAGING);
657     }
658
659     assert((((size_t) block) & (ALIGNMENT - 1)) == 0);
660
661     void* ptr = AddVoidPtr(block, METADATA_LEADING);
662     if (flags & HEAP_ZERO) {
663         inline_memset(ptr, 0, size);
664     }
665
666     return ptr;
667 }
668
669 void* AllocHeap(size_t size) {
670     return AllocHeapEx(size, HEAP_NO_FAULT);
671 }
672
673 void* AllocHeapZero(size_t size) {
674     return AllocHeapEx(size, HEAP_ZERO | HEAP_NO_FAULT);
675 }
676
677 void FreeHeap(void* ptr) {
678     MAX_IRQL(IRQL_SCHEDULER);
679
680     if (ptr == NULL) {
681         return;
682     }
683
684     struct block* block = SubVoidPtr(ptr, METADATA_LEADING);
685     bool pagable = IsOnSwappableHeap(block);
686     block->prev = NULL;
687     block->next = NULL;
688

```



```

689     bool acquired = LockHeap(pagable);
690
691     AddBlock(block);
692
693 #ifndef NDEBUG
694     outstanding_allocations--;
695 #endif
696
697     if (acquired) {
698         UnlockHeap(pagable);
699     }
700 }
701
702 void ReinitHeap(void) {
703     heap_lock = CreateMutex("heap");
704 }
705
706 void InitHeap(void) {
707     InitSpinlock(&heap_spinlock, "heapspin", IRQL_SCHEDULER);
708     InitSpinlock(&heap_locker_lock, "locker", IRQL_HIGH);
709 }

```

File: ./mem/physical.c

```

1
2 #include <arch.h>
3 #include <physical.h>
4 #include <diskcache.h>
5 #include <common.h>
6 #include <spinlock.h>
7 #include <assert.h>
8 #include <string.h>
9 #include <irql.h>
10 #include <log.h>
11 #include <virtual.h>
12 #include <panic.h>
13
14 static struct spinlock phys_lock;
15
16 /*
17  * One bit per page. Lower bits refer to lower pages. A clear bit indicates
18  * the page is unavailable (allocated / non-RAM), and a set bit indicates the
19  * page is free.
20  */
21 #define MAX_MEMORY_PAGES (ARCH_MAX_RAM_KBS / ARCH_PAGE_SIZE * 1024)
22 #define BITS_PER_ENTRY (sizeof(size_t) * 8)
23 #define BITMAP_ENTRIES (MAX_MEMORY_PAGES / BITS_PER_ENTRY)
24 static size_t allocation_bitmap[BITMAP_ENTRIES];
25
26 /*
27  * Stores pages that are available for us to allocate. If set to NULL, then we
28  * have yet to initialise the physical memory manager. The stack grows upward,
29  * and the pointer is incremented after writing the value on a push. The stack
30  * stores physical page numbers (indexes) instead of addresses.
31  */
32 static size_t* allocation_stack = NULL;
33 static size_t allocation_stack_pointer = 0;
34
35 /*
36  * Once we get below this number, we will start evicting pages.
37  */
38 #define NUM_EMERGENCY_PAGES 32
39
40 /*
41  * The number of physical pages available (free) remaining, and total, in the
42  * system.
43  */
44 static size_t pages_left = 0;
45 static size_t total_pages = 0;
46
47 /*
48  * The highest physical page number that exists on this system. Gets set during
49  * InitPhys() when scanning the system's memory map.
50  */
51 static size_t highest_valid_page_index = 0;
52
53 static inline bool IsBitmapEntryFree(size_t index) {
54     size_t base = index / BITS_PER_ENTRY;
55     size_t offset = index % BITS_PER_ENTRY;
56     return allocation_bitmap[base] & (1 << offset);
57 }
58
59 static inline void AllocateBitmapEntry(size_t index) {
60     assert(IsBitmapEntryFree(index));
61
62     size_t base = index / BITS_PER_ENTRY;
63     size_t offset = index % BITS_PER_ENTRY;
64     allocation_bitmap[base] |= (1 << offset);
65 }
66
67 static inline void DeallocateBitmapEntry(size_t index) {
68     assert(!IsBitmapEntryFree(index));
69
70     size_t base = index / BITS_PER_ENTRY;
71     size_t offset = index % BITS_PER_ENTRY;
72     allocation_bitmap[base] |= 1 << offset;
73 }
74
75 static inline void PushIndex(size_t index) {
76     assert(index <= highest_valid_page_index);
77     allocation_stack[allocation_stack_pointer++] = index;
78 }
79
80 static inline size_t PopIndex(void) {
81     assert(allocation_stack_pointer != 0);
82     return allocation_stack[--allocation_stack_pointer];
83 }
84
85 /*
86  * Removes an entry from the stack by value. Only to be used when absolutely
87  * required, as it has O(n) runtime and is therefore very slow.
88  */
89 static void RemoveStackEntry(size_t index) {
90     for (size_t i = 0; i < allocation_stack_pointer; ++i) {
91         if (allocation_stack[i] == index) {
92             memmove(allocation_stack + i, allocation_stack + i + 1, (--allocation_stack_pointer - i) * sizeof(size_t));
93             return;
94         }
95     }
96 }
97
98 /**
99  * Deallocates a page of physical memory that was allocated with AllocPhys().
100  * Does not affect virtual mappings - that should be taken care of before
101  * deallocating. Address must be page aligned.
102  */

```

```

103 void DeallocPhys(size_t addr) {
104     MAX_IQRL(IQRL_SCHEDULER);
105     assert(addr % ARCH_PAGE_SIZE == 0);
106
107     size_t page = addr / ARCH_PAGE_SIZE;
108
109     AcquireSpinlockIrql(&phys_lock);
110
111     ++pages_left;
112     DeallocateBitmapEntry(page);
113     if (allocation_stack != NULL) {
114         PushIndex(page);
115     }
116     ReleaseSpinlockIrql(&phys_lock);
117
118     if (pages_left > NUM_EMERGENCY_PAGES * 2) {
119         SetDiskCaches(DISKCACHE_NORMAL);
120     }
121 }
122
123 /**
124  * Deallocates a section of physical memory that was allocated with AllocPhysContinuous(). The entire block
125  * of memory must be deallocated at once, i.e. the start address of the memory should be passed in. Does not
126  * affect virtual mappings - that should be taken care of before deallocating.
127  *
128  * @param addr The address of the section of memory to deallocate. Must be page-aligned.
129  * @param size The size of the allocation. This should be the same value that was passed into AllocPhysContinuous().
130  */
131 void DeallocPhysContiguous(size_t addr, size_t bytes) {
132     for (size_t i = 0; i < BytesToPages(bytes); ++i) {
133         DeallocPhys(addr);
134         addr += ARCH_PAGE_SIZE;
135     }
136 }
137
138 static void EvictPagesIfNeeded(void*) {
139     EXACT_IQRL(IQRL_STANDARD);
140
141     extern int handling_page_fault;
142     if (handling_page_fault > 0) {
143         return;
144     }
145
146     if (pages_left < NUM_EMERGENCY_PAGES) {
147         SetDiskCaches(DISKCACHE_TOS);
148     } else if (pages_left < NUM_EMERGENCY_PAGES * 3 / 2) {
149         SetDiskCaches(DISKCACHE_REDUCE);
150     }
151
152     int timeout = 0;
153     while (pages_left < NUM_EMERGENCY_PAGES && timeout < 5) {
154         handling_page_fault++;
155         EvictVirt();
156         handling_page_fault--;
157         ++timeout;
158     }
159 }
160
161 size_t AllocPhys(void) {
162     MAX_IQRL(IQRL_SCHEDULER);
163
164     AcquireSpinlockIrql(&phys_lock);
165
166     if (pages_left == 0) {
167         Panic(PANIC_OUT_OF_PHYS);
168     }
169     if (pages_left <= NUM_EMERGENCY_PAGES) {
170         DeferUntilIrql(IQRL_STANDARD, EvictPagesIfNeeded, NULL);
171     }
172
173     size_t index = 0;
174     if (allocation_stack == NULL) {
175         /*
176          * No stack yet, so must use the bitmap. No point optimising this as
177          * only used during boot.
178          */
179         while (!IsBitmapEntryFree(index)) {
180             index = (index + 1) % MAX_MEMORY_PAGES;
181         }
182     } else {
183         index = PopIndex();
184     }
185
186     AllocateBitmapEntry(index);
187     --pages_left;
188     ReleaseSpinlockIrql(&phys_lock);
189
190     return index * ARCH_PAGE_SIZE;
191 }
192
193 /**
194  * Allocates a section of contiguous physical memory, that may or may not have
195  * requirements as to where the memory can be located. Must only be called after
196  * a call to ReinitPhys() is made. Deallocation should be done by
197  * DeallocPhysContiguous(), passing in the same size value as passed into
198  * AllocPhysContiguous() on allocation. Will not cause pages to be evicted from
199  * RAM, so sufficient memory must exist on the system for this to succeed.
200  *
201  * @param bytes The size of the allocation, in bytes.
202  * @param min_addr Allocated memory will not contain any addresses lower than
203  * this value.
204  * @param max_addr Allocated memory will not contain any addresses greater or
205  * equal to this value. For no maximum, set to 0.
206  * @param boundary Allocated memory will not contain any addresses that are an
207  * integer multiple of this value (although it may start at an
208  * integer multiple of this address). If there are no boundary
209  * requirements, set this to 0.
210  * @return The start address of the returned physical memory area. If the
211  * request could not be satisfied, 0 is returned.
212  */
213 size_t AllocPhysContiguous(size_t bytes, size_t min_addr, size_t max_addr, size_t boundary) {
214     /*
215      * This function should not be called before the stack allocator is setup.
216      * (There is no need for InitVirt() to use this function, and so checking
217      * here removes a check that would have to be done in a loop later).
218      */
219     if (allocation_stack == NULL) {
220         return 0;
221     }
222
223     size_t pages = BytesToPages(bytes);
224     size_t min_index = (min_addr + ARCH_PAGE_SIZE - 1) / ARCH_PAGE_SIZE;
225     size_t max_index = max_addr == 0 ? highest_valid_page_index + 1 : max_addr / ARCH_PAGE_SIZE;
226     size_t count = 0;
227
228     AcquireSpinlockIrql(&phys_lock);
229
230     /*
231      * We need to check we won't try to over-allocate memory, or allocate so
232

```

```

233     * much memory that it puts us in a critical position.
234     */
235     if (pages + NUM_EMERGENCY_PAGES >= pages_left) {
236         ReleaseSpinlockIrql(&phys_lock);
237         return 0;
238     }
239
240     for (size_t index = min_index; index < max_index; ++index) {
241         /*
242          * Reset the counter if we are no longer contiguous, or if we have cross
243          * a boundary that we can't cross.
244          */
245         if (!IsBitmapEntryFree(index) || (boundary != 0 && (index % (boundary / ARCH_PAGE_SIZE) == 0))) {
246             count = 0;
247             continue;
248         }
249
250         ++count;
251         if (count == pages) {
252             /*
253              * Go back to the start of the section and mark it all as allocated.
254              */
255             size_t start_index = index - count + 1;
256             while (start_index <= index) {
257                 AllocateBitmapEntry(start_index);
258                 RemoveStackEntry(start_index);
259                 ++start_index;
260             }
261
262             ReleaseSpinlockIrql(&phys_lock);
263             return start_index * ARCH_PAGE_SIZE;
264         }
265     }
266
267     ReleaseSpinlockIrql(&phys_lock);
268     return 0;
269 }
270
271 /**
272  * Initialises the physical memory manager for the first time. Must be called
273  * before any other memory management function is called. It determines what
274  * memory is available on the system and prepares the O(n) bitmap allocator.
275  * This will be slow, but is only needed until ReinitHeap() gets called.
276  */
277 void InitPhys(void) {
278     InitSpinlock(&phys_lock, "phys", IRQL_SCHEDULER);
279
280     /*
281     * Scan the memory tables and fill in the memory that is there.
282     */
283     while (true) {
284         struct arch_memory_range* range = ArchGetMemory();
285
286         if (range == NULL) {
287             /* No more memory exists */
288             break;
289         }
290     } else {
291         /*
292          * Must round the start address up so we don't include memory outside
293          * the region.
294          */
295         size_t first_page = (range->start + ARCH_PAGE_SIZE - 1) / ARCH_PAGE_SIZE;
296         size_t last_page = (range->start + range->length) / ARCH_PAGE_SIZE;
297
298         while (first_page < last_page && first_page < MAX_MEMORY_PAGES) {
299             DeallocateBitmapEntry(first_page);
300             ++pages_left;
301             ++total_pages;
302
303             if (first_page > highest_valid_page_index) {
304                 highest_valid_page_index = first_page;
305             }
306
307             ++first_page;
308         }
309     }
310 }
311
312 static void ReclaimBitmapSpace(void) {
313     /*
314     * Save extra physical memory on by deallocating the memory in the bitmap
315     * that can't be reached (due to the system not having memory that goes up
316     * that high).
317     */
318     size_t unreachable_pages = MAX_MEMORY_PAGES - (highest_valid_page_index + 1);
319     size_t unreachable_entries = unreachable_pages / BITS_PER_ENTRY;
320     size_t unreachable_bitmap_pages = unreachable_entries / ARCH_PAGE_SIZE;
321
322     size_t end_bitmap = ((size_t) allocation_bitmap) + sizeof(allocation_bitmap);
323
324     /*
325     * Round down, otherwise other kernel data in the same page as the end of
326     * the bitmap will also be counted as 'free', causing memory corruption.
327     */
328     size_t unreachable_region = ((end_bitmap - ARCH_PAGE_SIZE * unreachable_bitmap_pages) & ~(ARCH_PAGE_SIZE - 1));
329
330     while (num_unreachable_bitmap_pages-- > 0) {
331         DeallocPhys(ArchVirtualToPhysical(unreachable_region));
332         unreachable_region += ARCH_PAGE_SIZE;
333         ++total_pages;
334     }
335 }
336
337 /**
338  * Reinitialises the physical memory manager with a constant-time page allocation system.
339  * Must be called after virtual memory has been initialised. Must only be called once. Must
340  * be called before calling AllocPageContiguous() is called. Should be called as soon as
341  * possible after virtual memory is available.
342  */
343 void ReinitPhys(void) {
344     assert(allocation_stack == NULL);
345
346     allocation_stack = (size_t*) MapVirt(0, 0, (highest_valid_page_index + 1) * sizeof(size_t), VM_READ | VM_WRITE | VM_LOCK, NULL, 0);
347
348     for (size_t i = 0; i < MAX_MEMORY_PAGES; ++i) {
349         if (IsBitmapEntryFree(i)) {
350             PushIndex(i);
351         }
352     }
353
354     ReclaimBitmapSpace();
355 }
356
357 size_t GetTotalPhysKilobytes(void) {
358     return total_pages * ARCH_PAGE_SIZE / 1024;
359 }
360
361 size_t GetFreePhysKilobytes(void) {

```

```

363     return pages_left * (ARCH_PAGE_SIZE / 1024);
364 }

```

File: ./mem/virtual.c

```

1
2  /**
3   * mem/virtual.c - Virtual Memory Manager
4   */
5
6  #include <virtual.h>
7  #include <avl.h>
8  #include <heap.h>
9  #include <common.h>
10 #include <arch.h>
11 #include <physical.h>
12 #include <assert.h>
13 #include <panic.h>
14 #include <string.h>
15 #include <timer.h>
16 #include <driver.h>
17 #include <thread.h>
18 #include <spinlock.h>
19 #include <log.h>
20 #include <sys/types.h>
21 #include <irq.h>
22 #include <irq.h>
23 #include <cpu.h>
24 #include <transfer.h>
25 #include <console.h>
26 #include <swapfile.h>
27 #include <vfs.h>
28 #include <errno.h>
29
30 // TODO: lots of locks! especially the global cpu one
31
32 static size_t SplitLargePageEntryIntoMultiple(struct vas* vas, size_t virtual, struct vas_entry* entry, int num_to_leave);
33 static struct vas_entry* GetVirtEntry(struct vas* vas, size_t virtual);
34
35 /**
36 * Stores a pointer to any kernel VAS. Ensures that when processes are destroyed, we are using a VAS
37 * that is different from the VAS that's being deleted.
38 */
39 static struct vas* kernel_vas;
40
41 /**
42 * Used in debugging as to print out the contents of the mappings tree.
43 */
44 void AvlPrinter(void* entry_) {
45     struct vas_entry* entry = (struct vas_entry*) entry_;
46
47     LogWriteSerial(
48         "[v: 0x%X, p: 0x%X; acrl: %d&d&d&d. ref: %d]; ",
49         entry->virtual, entry->physical, entry->allocated, entry->cow, entry->in_ram, entry->lock, entry->ref_count
50     );
51 }
52
53 /**
54 * Whether or not virtual memory is available for use. Can be read with IsVirtInitialised(), and is set when
55 * InitVirt() has completed.
56 *
57 * Does not have a lock, as only the bootstrap CPU should be modifying it, and this happens before threads
58 * are set up. Once set to true, it is never changed again, so there is no read/write problems.
59 */
60 static bool virt_initialised = false;
61
62 /**
63 * Used by the mappings tree as its comparison operator. Allows us to properly maintain the AVL properties,
64 * and for 'GetVirtEntry' to work.
65 *
66 * Both parameters should be cast to 'struct vas_entry*' and then used as normal.
67 */
68 static int VirtAvlComparator(void* a, void* b) {
69     struct vas_entry* a_entry = (struct vas_entry*) a;
70     struct vas_entry* b_entry = (struct vas_entry*) b;
71
72     assert((a_entry->virtual & (ARCH_PAGE_SIZE - 1)) == 0);
73     assert((b_entry->virtual & (ARCH_PAGE_SIZE - 1)) == 0);
74
75     /*
76      * Check for overlapping regions for multi-mapping entries, and count that as equal. This allows us
77      * to return the correct entry if one of them is part of a multi-mapping entry.
78      */
79     size_t a_page = a_entry->virtual / ARCH_PAGE_SIZE;
80     size_t b_page = b_entry->virtual / ARCH_PAGE_SIZE;
81     if (a_page == b_page && a_page < b_page + b_entry->num_pages) {
82         return 0;
83     }
84     if (b_page >= a_page && b_page < a_page + a_entry->num_pages) {
85         return 0;
86     }
87
88     return COMPARE_SIGN(a_entry->virtual, b_entry->virtual);
89 }
90
91 /**
92 * Initialises a virtual address space in an already allocated section of memory.
93 *
94 * @param vas The memory to initialise a virtual address space object in.
95 * @param flags Can be 0 or VAS_NO_ARCH_INIT. If VAS_NO_ARCH_INIT is provided, then no architecture-specific
96 *             code will be called. This flag should only be set if called by architecture-specific functions
97 *             (e.g. to create the initial address space). Normally, 0 should be passed in.
98 *
99 * @maxirq Irql_SCHEDULER
100 */
101 void CreateVasEx(struct vas* vas, int flags) {
102     MAX_IRQL(IRQL_SCHEDULER);
103
104     vas->mappings = AvlTreeCreate();
105
106     /*
107      * We are in for a world of hurt if someone is able to page fault while
108      * holding the lock on a virtual address space, so better make it IRQL_SCHEDULER.
109      */
110     InitSpinlock(&vas->lock, "vas", IRQL_SCHEDULER);
111     AvlTreeSetComparator(vas->mappings, VirtAvlComparator);
112     if (!(flags & VAS_NO_ARCH_INIT)) {
113         ArchInitVas(vas);
114     }
115 }
116
117 /**
118 * Allocates and initialises a new virtual address space.
119 *
120 * @return The virtual address space which was created.
121 */

```

```

122 * @maxirq Irql_SCHEDULER
123 */
124 struct vas* CreateVas() {
125     MAX_IRQL(IRQL_SCHEDULER);
126
127     struct vas* vas = AllocHeap(sizeof(struct vas));
128     CreateVasEx(vas, 0);
129     return vas;
130 }
131
132 struct defer_disk_access {
133     struct open_file* file;
134     struct vas_entry* entry;
135     off_t offset;
136     size_t address;
137     int direction;
138     bool deallocate_swap_on_read;
139 };
140
141 static void PerformDeferredAccess(void* data) {
142     // TODO: see comment in BringIntoMemoryFromFile
143
144     struct defer_disk_access* access = (struct defer_disk_access*) data;
145
146     bool write = access->direction == TRANSFER_WRITE;
147     size_t target_address = access->address;
148     if (!write) {
149         /*
150          * If we're reading, the page is not yet allocated or in memory (this is so we don't have other threads
151          * trying to use the partially-filled page). Therefore, we allocate a temporary page to write the data in, and
152          * then we can allocate the page and copy the data while we hold a lock.
153          *
154          * We can't just allocate the proper page entry now, as we can't hold the spinlock over the call to ReadFile.
155          */
156         target_address = MapVirt(0, 0, ARCH_PAGE_SIZE, VM_LOCK | VM_READ | VM_WRITE, NULL, 0);
157     }
158
159     struct transfer tr = CreateKernelTransfer((void*) target_address, ARCH_PAGE_SIZE, access->offset, access->direction);
160
161     int res = write ? WriteFile : ReadFile(access->file, &tr);
162     if (res != 0) {
163         /*
164          * TODO: it's not actually always a failure. the only 'panic' condition is when it involves
165          * the swapfile, but this code is also used for dealing with normal file-mapped pages.
166          *
167          * for file-mapped pages, failures due to reading past the end of the file should always
168          * be okay - we need to fill the rest of the page with zero though (even if that page has
169          * no file data on it, e.g. if we read really past the end of the array).
170          */
171         if (access->entry->swapfile) {
172             Panic(PANIC_DISK_FAILURE_ON_SWAPFILE);
173         } else {
174             // I think for reads, it's okay to not do anything here on error, and just make use of
175             // the number of bytes that were actually transferred (and therefore complete failure means
176             // we just end up with a blanked-out page being allocated).
177             Panic(PANIC_NOT_IMPLEMENTED);
178         }
179     }
180
181     if (write) {
182         UnmapVirt(access->address, ARCH_PAGE_SIZE);
183     } else {
184         /*
185          * Now we can actually lock the page and allocate the actual mapping.
186          */
187         struct vas* vas = GetVas();
188         AcquireSpinlockIrql(&vas->lock);
189
190         struct vas_entry* entry = GetVirtEntry(vas, access->address);
191         assert(entry->num_pages == 1);
192         assert(entry->swapfile || entry->file);
193
194         entry->lock = true;
195         entry->physical = AllocPhys();
196         entry->allocated = true;
197         entry->allow_temp_write = true;
198         entry->in_ram = true;
199         entry->swapfile = false;
200         ArchUpdateMapping(vas, entry);
201         ArchFlushTlb(vas);
202
203         // TODO: this should use the actual amount that was read...
204
205         inline_memcpy((void*) access->address, (const char*) target_address, ARCH_PAGE_SIZE);
206
207         entry->allow_temp_write = false;
208
209         /*
210          * If it was on the swapfile, we now need to mark that slot in the swapfile as free for future use.
211          */
212         if (access->deallocate_swap_on_read) {
213             DeallocateSwapfileIndex(access->offset / ARCH_PAGE_SIZE);
214         }
215
216         /*
217          * Don't perform relocations on the first load, as the first load will be when 'proper' relocation
218          * happens (i.e. the 'all at once' relocations) - and therefore the quick relocation table will not
219          * be created yet and we'll crash.
220          *
221          * The reason we can't just not do the initial big relocation and make it all work though demand loading
222          * is because not all pages with driver code/data end up being marked as VM_RELOCATABLE (e.g. for small
223          * parts of data segments, etc.).
224          */
225         bool needs_relocations = entry->relocatable && !entry->first_load;
226
227         ArchUpdateMapping(vas, entry);
228         ArchFlushTlb(vas);
229
230         /*
231          * Need to keep page locked if we're doing relocations on it - otherwise by the time that we actually
232          * load in all the data we need to do the relocations (e.g. ELF headers, the symbol table), we have probably
233          * already swapped out the page we are relocating (which leads to us getting nowhere).
234          */
235         if (!needs_relocations) {
236             entry->first_load = false;
237             entry->load_in_progress = false;
238             entry->lock = false;
239         }
240         ReleaseSpinlockIrql(&vas->lock);
241
242         UnmapVirt(target_address, ARCH_PAGE_SIZE);
243
244         if (needs_relocations) {
245             LogWriteSerial(" ----> ABOUT TO PERFORM RELOCATION FIXUPS\n");
246             PerformDriverRelocationOnPage(vas, entry->relocation_base, access->address);
247             AcquireSpinlockIrql(&vas->lock);
248             entry->first_load = false;
249             entry->load_in_progress = false;
250             UnlockVirtEx(vas, access->address);
251         }

```

```

252     ReleaseSpinlockIrql(&vas->lock);
253     LogWriteSerial(" ----> PERFORMED RELOCATION FIXUPS\n");
254 }
255
256     LogWriteSerial(" ----> FINISHED RELOADING FROM DISK 0x%X\n", entry->virtual);
257 }
258
259     FreeHeap(access);
260 }
261
262 /**
263  * Given a virtual page, it defers a write to disk. It creates a copy of the virtual page, so that it may be safely
264  * deleted as soon as this gets called.
265  */
266 static void DeferDiskWrite(size_t old_addr, struct open_file* file, off_t offset) {
267     size_t new_addr = MapVirt(0, 0, ARCH_PAGE_SIZE, VM_LOCK | VM_READ | VM_WRITE | VM_RECURSIVE, NULL, 0);
268     inline_memcpy((void*) new_addr, (const char*) old_addr, ARCH_PAGE_SIZE);
269
270     struct defer_disk_access* access = AllocHeap(sizeof(struct defer_disk_access));
271     access->address = new_addr;
272     access->file = file;
273     access->direction = TRANSFER_WRITE;
274     access->offset = offset;
275     access->deallocate_swap_on_read = false;
276     DeferUntilIrql(IRQL_STANDARD_HIGH_PRIORITY, PerformDeferredAccess, (void*) access);
277 }
278
279 static void DeferDiskRead(size_t new_addr, struct open_file* file, off_t offset, bool deallocate_swap_on_read) {
280     struct defer_disk_access* access = AllocHeap(sizeof(struct defer_disk_access));
281     access->address = new_addr;
282     access->file = file;
283     access->direction = TRANSFER_READ;
284     access->offset = offset;
285     access->deallocate_swap_on_read = deallocate_swap_on_read;
286     DeferUntilIrql(IRQL_STANDARD_HIGH_PRIORITY, PerformDeferredAccess, (void*) access);
287 }
288
289 /**
290  * Evicts a particular page mapping from virtual memory, freeing up its physical page (if it had one).
291  * This will often involve accessing the disk to put it on swapfile (or save modifications to a file-backed
292  * page).
293  *
294  * @param vas The virtual address space that we're evicting from. Does not have to be the current one.
295  * @param entry The virtual page to remove from virtual memory.
296  *
297  * @maxirql IRQL_SCHEDULER
298  */
299 void EvictPage(struct vas* vas, struct vas_entry* entry) {
300     MAX_IRQL(IRQL_PAGE_FAULT);
301
302     assert(!entry->lock);
303     assert(!entry->cow);
304
305     AcquireSpinlockIrql(&vas->lock);
306
307     if (entry->in_ram) {
308         /*
309          * Nothing happens, as this page isn't even in RAM.
310          */
311     } else if (entry->file) {
312         /*
313          * We will just reload it from disk next time.
314          */
315     }
316
317     if (entry->write && entry->relocatable) {
318         DeferDiskWrite(entry->virtual, entry->file_node, entry->file_offset);
319     }
320
321     entry->in_ram = false;
322     entry->allocated = false;
323     DeallocPhys(entry->physical);
324     ArchUnmap(vas, entry);
325     ArchFlushTlb(vas);
326
327 } else {
328     /*
329      * Otherwise, we need to mark it as swapfile.
330      */
331     entry->in_ram = false;
332     entry->swapfile = true;
333     entry->allocated = false;
334
335     uint64_t offset = AllocateSwapfileIndex() * ARCH_PAGE_SIZE;
336
337     //PutsConsole("PAGE OUT\n");
338     LogWriteSerial(" ----> WRITING VIRT 0x%X TO SWAP: DISK INDEX 0x%X (offset 0x%X)\n", entry->virtual, (int) offset / ARCH_PAGE_SIZE, (int) offset);
339     DeferDiskWrite(entry->virtual, GetSwapfile(), offset);
340     entry->swapfile_offset = offset;
341
342     ArchUnmap(vas, entry);
343     DeallocPhys(entry->physical);
344     ArchFlushTlb(vas);
345 }
346
347     ReleaseSpinlockIrql(&vas->lock);
348 }
349
350 /**
351  * Lower value means it should be swapped out first.
352  */
353 static int GetPageEvictionRank(struct vas* vas, struct vas_entry* entry) {
354     (void) vas;
355
356     /*
357      * Want to evict in this order:
358      * - file and non-writable
359      * - file and writable
360      * - non-writable
361      * - writable
362      *
363      * When we have a way of dealing with accessed / dirty, it should be in this order:
364      *
365      * 0 VM_EVICT_FIRST
366      * 10 FILE, NON-WRITABLE, NON-ACCESSED
367      * 20 FILE, WRITABLE, NON-DIRTY, NON-ACCESSED
368      * 30 FILE, NON-WRITABLE, ACCESSED
369      * 40 FILE, WRITABLE, NON-DIRTY, ACCESSED
370      * 50 NORMAL, NON-DIRTY, NON-ACCESSED
371      * 60 NORMAL, NON-DIRTY, ACCESSED
372      * 70 FILE, WRITABLE, DIRTY
373      * 80 NORMAL, DIRTY
374      * 90 COW
375      * 150 RELOCATABLE
376      *
377      * Globals add 3 points.
378      */
379
380     bool accessed;
381     bool dirty;

```

```

382 ArchGetPageUsageBits(vas, entry, &accessed, &dirty);
383 ArchSetPageUsageBits(vas, entry, false, false);
384
385 int penalty = (entry->global ? 3 : 0) + entry->times_swapped * 8;
386
387 if (entry->evict_first) {
388     return entry->times_swapped;
389 }
390 } else if (entry->relocatable) {
391     return 150;
392 }
393 } else if (entry->cow) {
394     return 90 + penalty;
395 }
396 } else if (entry->file && (entry->write)) {
397     return (accessed ? 30 : 10) + penalty;
398 }
399 } else if (entry->file && entry->write) {
400     return (dirty ? 70 : (accessed ? 40 : 20)) + penalty;
401 }
402 } else if (!dirty) {
403     return (accessed ? 60 : 50) + penalty;
404 }
405 } else {
406     return 80 + penalty;
407 }
408 }
409
410 struct eviction_candidate {
411     struct vas* vas;
412     struct vas_entry* entry;
413 };
414
415 #define PREV_SWAP_LIMIT 24
416
417 void FindVirtToEvictFromSubtree(struct vas* vas, struct avl_node* node, int* lowest_rank, struct eviction_candidate* lowest_ranked, int* count, struct vas_entry
418     static uint8_t rand = 0;
419
420     if (node == NULL) {
421         return;
422     }
423
424     if (*lowest_rank < 10) {
425         /*
426          * No need to look anymore - we've already a best possible page.
427          */
428         return;
429     }
430
431     *count += 1;
432
433     /*
434      * After scanning through 500 entries, we'll allow early exits for less optimal pages.
435      */
436     int limit = (((*count - 500) / 75) + 10);
437     if (*count > 500 && *lowest_rank < limit) {
438         return;
439     }
440
441     struct vas_entry* entry = AvlTreeGetData(node);
442     if ((entry->lock && entry->allocated) {
443         int rank = GetPageEvictionRank(vas, entry);
444
445         /*
446          * To ensure we mix up who gets evicted, when there's an equality, we use it 1/4 times.
447          * It is likely there are more than 4 to replace, so this ensures that we cycle through many of them.
448          */
449         bool equal = rank == *lowest_rank;
450         if (equal) {
451             equal = (rand++ % 3) == 0;
452         }
453
454         bool prev_swap = false;
455         for (int i = 0; i < PREV_SWAP_LIMIT; ++i) {
456             if (prev_swaps[i] == entry) {
457                 prev_swap = true;
458                 break;
459             }
460         }
461
462         if ((rank < *lowest_rank || equal) && !prev_swap) {
463             lowest_ranked->vas = vas;
464             lowest_ranked->entry = entry;
465             *lowest_rank = rank;
466
467             if (rank == 0) {
468                 return;
469             }
470         }
471     }
472
473     FindVirtToEvictFromSubtree(vas, AvlTreeGetLeft(node), lowest_rank, lowest_ranked, count, prev_swaps);
474     FindVirtToEvictFromSubtree(vas, AvlTreeGetRight(node), lowest_rank, lowest_ranked, count, prev_swaps);
475 }
476
477 void FindVirtToEvictFromAddressSpace(struct vas* vas, int* lowest_rank, struct eviction_candidate* lowest_ranked, bool include_globals, struct vas_entry** prev_
478     int count = 0;
479     FindVirtToEvictFromSubtree(vas, AvlTreeGetRootNode(vas->mappings), lowest_rank, lowest_ranked, &count, prev_swaps);
480     if (include_globals)
481         FindVirtToEvictFromSubtree(vas, AvlTreeGetRootNode(GetCpu()->global_vas_mappings), lowest_rank, lowest_ranked, &count, prev_swaps);
482 }
483
484 /**
485  * Searches through virtual memory (that doesn't necessarily have to be in the current virtual address space),
486  * and finds and evicts a page of virtual memory, to try free up physical memory.
487  *
488  * @maxirql IRQL_STANDARD
489  */
490
491 void EvictVirt(void) {
492     MAX_IRQL(IRQL_PAGE_FAULT);
493
494     if (GetSwapfile() == NULL) {
495         return;
496     }
497
498     // TODO: we need to ensure that EvictVirt(), when called from the defer, does not evict any pages that were just
499     // loaded in!! This is an issue when we need to perform relocations during page faults, as that brings in a
500     // whole heap of other pages, and that often causes TryEvictPages() to straight away get rid of the page we just
501     // loaded in. Alternatively, TryEvictPages() can be a NOP the first time it is called after a page fault.
502     // This would give the code on the page that we loaded in time to 'progress' before being swapped out again.
503
504     /*
505     void TryEvictPages() {
506         if (had_page_fault) {
507             had_page_fault = false;
508             return;
509         }
510         EvictVirt()
511     }

```

```

512
513 void HandlePageFault() {
514     had_page_fault = true;
515 }
516
517 /*
518 // don't allow any of the last 8 swaps to be repeated (as an instruction may require at least 6 pages on x86
519 // if it straddles many boundaries)
520 static struct vas_entry* previous_swaps[PREV_SWAP_LIMIT] = {0};
521 static int swap_num = 0;
522
523 int lowest_rank = 10000;
524 struct eviction_candidate lowest_ranked;
525 lowest_ranked.entry = NULL;
526
527 AcquireSpinlockIrql(&GetVas()->lock);
528 FindVirtToEvictFromAddressSpace(GetVas(), &lowest_rank, &lowest_ranked, true, previous_swaps);
529 ReleaseSpinlockIrql(&GetVas()->lock);
530
531 // TODO: go through other address spaces
532
533 while (false) {
534     struct vas* vas = NULL;
535     if (vas != GetVas())
536         FindVirtToEvictFromAddressSpace(GetVas(), &lowest_rank, &lowest_ranked, false, previous_swaps);
537 }
538
539 if (lowest_ranked.entry != NULL) {
540     previous_swaps[swap_num++ % PREV_SWAP_LIMIT] = lowest_ranked.entry;
541     EvictPage(lowest_ranked.vas, lowest_ranked.entry);
542     lowest_ranked.entry->times_swapped++;
543 }
544 }
545
546 static void InsertIntoAvl(struct vas* vas, struct vas_entry* entry) {
547     assert(IsSpinlockHeld(&vas->lock));
548
549     if (entry->global) {
550         AcquireSpinlockIrql(&GetCpu()->global_mappings_lock);
551         AvlTreeInsert(GetCpu()->global_vas_mappings, entry);
552         ReleaseSpinlockIrql(&GetCpu()->global_mappings_lock);
553     } else {
554         AvlTreeInsert(vas->mappings, entry);
555     }
556 }
557
558 static void DeleteFromAvl(struct vas* vas, struct vas_entry* entry) {
559     assert(IsSpinlockHeld(&vas->lock));
560     if (entry->global) {
561         AcquireSpinlockIrql(&GetCpu()->global_mappings_lock);
562         AvlTreeDelete(GetCpu()->global_vas_mappings, entry);
563         ReleaseSpinlockIrql(&GetCpu()->global_mappings_lock);
564     } else {
565         AvlTreeDelete(vas->mappings, entry);
566     }
567 }
568
569 /**
570 * Adds a virtual page mapping to the specified virtual address space. This will add it both to the mapping tree
571 * and the architectural paging structures (so that page faults can be raised, etc., if there is no backing yet).
572 *
573 * @param vas The virtual address space to map this page to
574 * @param physical Only used if VM_LOCK is specified in flags. Determines the physical page that will back the
575 * virtual mapping. If VM_LOCK is set, and this is 0, then a physical page will be allocated. If
576 * VM_LOCK is set, and this is non-zero, then that physical address will be used.
577 * @param virtual The virtual address to map the memory to. This should be non-zero.
578 * @param flags Various bitflags to affect the attributes of the mapping. Flags that are used here are:
579 * VM_READ : if set, the page will be marked as readable
580 * VM_WRITE : if set, the page will be marked as writable
581 * VM_USER : if set, then usermode can access this page without faulting
582 * VM_EXEC : if set, then code can be executed in this page
583 * VM_LOCK : if set, the page will immediately get a physical memory backing, and will not be
584 * paged out
585 * VM_FILE : if set, this page is file-backed
586 * All other flags are ignored by this function.
587 * @param file If VM_FILE is set, then the page is backed by this file, starting at the position specified by pos.
588 * @param pos If VM_FILE is set, then this is the offset into the file where the page is mapped to.
589 *
590 * @maxirql IRQL_SCHEDULER
591 */
592 static void AddMapping(struct vas* vas, size_t physical, size_t virtual, int flags, struct open_file* file, off_t pos, size_t number) {
593     MAX_IRQL(IRQL_SCHEDULER);
594
595     assert(!((file != NULL && (flags & VM_FILE) == 0)));
596
597     struct vas_entry* entry = AllocHeapZero(sizeof(struct vas_entry));
598     entry->allocated = false;
599
600     bool lock = flags & VM_LOCK;
601     entry->lock = lock;
602     entry->in_ram = lock;
603
604     size_t relocation_base = 0;
605     if (flags & VM_RELOCATABLE) {
606         relocation_base = physical;
607         physical = 0;
608     }
609
610     if (lock) {
611         /*
612          * We are not allowed to check if the physical page is allocated/free, because it might come
613          * from a VM_MAP_HARDWARE request, which can map non-RAM pages.
614          */
615         if (physical == 0) {
616             physical = AllocPhys();
617             entry->allocated = true;
618         }
619     }
620
621     /*
622      * MapVirt checks for conflicting flags and returns, so this code doesn't need to worry about that.
623      */
624     entry->virtual = virtual;
625     entry->times_swapped = 0;
626     entry->read = (flags & VM_READ) ? 1 : 0;
627     entry->write = (flags & VM_WRITE) ? 1 : 0;
628     entry->exec = (flags & VM_EXEC) ? 1 : 0;
629     entry->file = (flags & VM_FILE) ? 1 : 0;
630     entry->user = (flags & VM_USER) ? 1 : 0;
631     entry->evict_first = (flags & VM_EVICT_FIRST) ? 1 : 0;
632     entry->relocatable = (flags & VM_RELOCATABLE) ? 1 : 0;
633     entry->allow_temp_write = false;
634     entry->load_in_progress = false;
635     entry->global = !(flags & VM_LOCAL);
636     entry->physical = physical;
637     entry->ref_count = 1;
638     entry->file_offset = pos;
639 }

```



```

642     entry->file_node = file;
643     entry->swapfile = false;
644     entry->first_load = entry->relocatable;
645     entry->num_pages = number;
646
647     if (entry->relocatable) {
648         entry->relocation_base = relocation_base;
649     } else
650         entry->swapfile_offset = 0xDEADDEAD;
651 }
652
653 LogWriteSerial("Adding mapping at 0x%X to vas 0x%X - num is %d. flags = 0x%X, rxgu'lfia = %d%d%d%d%d' %d%d%d%d. p 0x%X.\n",
654     entry->virtual, vas, number, flags,
655     entry->read, entry->write, entry->exec, entry->global, entry->user, entry->lock, entry->file, entry->in_ram, entry->allocated,
656     entry->physical
657 );
658
659 /*
660  * TODO: later on, check if shared, and add phys->virt entry if needed
661  */
662
663 if ((flags & VM_RECURSIVE) == 0) {
664     AcquireSpinlockIrql(&vas->lock);
665 }
666
667 InsertIntoAvl(vas, entry);
668 ArchAddMapping(vas, entry);
669
670 if (entry->lock && (flags & VM_MAP_HARDWARE) == 0) {
671     if (GetVas() == vas) {
672         /*
673          * Need to zero out the page - this must happen on first load in, and as we have to load in
674          * locked pages now, we must do it now.
675          */
676         memset((void*) entry->virtual, 0, entry->num_pages * ARCH_PAGE_SIZE);
677     } else {
678         LogDeveloperWarning("yuck. PAGE HAS NOT BEEN ZEROED!\n");
679     }
680 }
681
682 if ((flags & VM_RECURSIVE) == 0) {
683     ReleaseSpinlockIrql(&vas->lock);
684 }
685
686 static bool IsRangeInUse(struct vas* vas, size_t virtual, size_t pages) {
687     bool in_use = false;
688
689     struct vas_entry dummy = {.num_pages = 1, .virtual = virtual};
690
691     /*
692      * We have to loop over the local one, and if it isn't there, the global one. We do this
693      * in separate loops to prevent the need to acquire both spinlocks at once, which could lead
694      * to a deadlock.
695      */
696
697     AcquireSpinlockIrql(&vas->lock);
698     for (size_t i = 0; i < pages; ++i) {
699         if (AvlTreeContains(vas->mappings, (void*) &dummy)) {
700             in_use = true;
701             break;
702         }
703         dummy.virtual += ARCH_PAGE_SIZE;
704     }
705     ReleaseSpinlockIrql(&vas->lock);
706
707     if (in_use) {
708         return true;
709     }
710
711     AcquireSpinlockIrql(&GetCpu()->global_mappings_lock);
712     dummy.virtual = virtual;
713     for (size_t i = 0; i < pages; ++i) {
714         if (AvlTreeContains(GetCpu()->global_vas_mappings, (void*) &dummy)) {
715             in_use = true;
716             break;
717         }
718         dummy.virtual += ARCH_PAGE_SIZE;
719     }
720     ReleaseSpinlockIrql(&GetCpu()->global_mappings_lock);
721
722     return in_use;
723 }
724
725 static size_t AllocVirtRange(struct vas* vas, size_t pages, int flags) {
726     /*
727      * TODO: make this deallocatable, and not x86 specific (with that memory address)
728      */
729     if (flags & VM_LOCAL) {
730         /*
731          * Also needs to use the vas to work out what's allocated in that vas
732          */
733         (void) vas;
734         static size_t hideous_allocator = 0x200000000U;
735         size_t retv = hideous_allocator;
736         hideous_allocator += pages * ARCH_PAGE_SIZE;
737         return retv;
738     } else {
739         /*
740          * TODO: this probably needs a global lock of some sort.
741          */
742         static size_t hideous_allocator = ARCH_KRNL_SBRK_BASE;
743         size_t retv = hideous_allocator;
744         hideous_allocator += pages * ARCH_PAGE_SIZE;
745         return retv;
746     }
747 }
748
749 static void FreeVirtRange(struct vas* vas, size_t virtual, size_t pages) {
750     (void) virtual;
751     (void) vas;
752     (void) pages;
753 }
754
755 /**
756  * Creates a virtual memory mapping.
757  *
758  * All mapped pages will be zeroed out (either on first use, or if locked, when allocated) - except if VM_MAP_HARDWARE or
759  * VM_FILE is set. If VM_FILE is set, reading beyond the end of the file, but within the page limit, will read zeroes.
760  *
761  * @param vas The virtual address space to map this page to
762  * @param physical Only used if VM_LOCK is specified in flags. Determines the physical page that will back the
763  * virtual mapping. If VM_LOCK is set, and this is 0, then a physical page will be allocated. If
764  * VM_LOCK is set, and this is non-zero, then that physical address will be used. In this instance,
765  * VM_MAP_HARDWARE must also be set. If VM_MAP_HARDWARE is not set, this value must be 0.
766  * @param virtual The virtual address to map the memory to. If this is 0, then a virtual memory region of the correct
767  * size will be allocated.
768  * @param pages The number of contiguous pages to map in this way
769  * @param flags Various bitflags to affect the attributes of the mapping. Flags that are used here are:
770  * VM_READ : if set, the page will be marked as readable
771  */

```

```

772 *          VM_WRITE      : if set, the page will be marked as writable. On some architectures, this may have the
773 *                          effect of implying VM_READ as well.
774 *          VM_USER       : if set, then usermode can access this page without faulting
775 *          VM_EXEC       : if set, then code can be executed in this page
776 *          VM_LOCK       : if set, the page will immediately get a physical memory backing, and will not be
777 *                          paged out.
778 *          VM_FILE       : if set, this page is file-backed. Cannot be combined with VM_MAP_HARDWARE.
779 *                          Cannot be combined with VM_LOCK.
780 *          VM_MAP_HARDWARE : if set, a physical address can be specified for the backing. If this flag is set,
781 *                          then VM_LOCK must also be set, and VM_FILE must be clear.
782 *          VM_LOCAL      : if set, it is only mapped into the current virtual address space. If set, it is mapped
783 *                          into the kernel virtual address space.
784 *          VM_RECURSIVE  : must be set if and only if this call to MapVirtEx is being called with the virtual
785 *                          address space lock already held. Does not affect the page, only the call to MapVirtEx
786 *                          When set, the lock is not automatically acquired or released as it is assumed to be
787 *                          already held.
788 *          VM_RELOCATABLE : if set, then this page will have driver relocations applied to it when it is swapped
789 *                          in. VM_FILE must be set as well. 'file' should be set to the driver's file.
790 *          VM_FIXED_VIRT : if set, then the virtual address specified in 'virtual' will be required to be used -
791 *                          if any page required for the mapping of this size is already allocated, the allocation
792 *                          will fail. If clear, then another virtual address may be used in order to satisfy a
793 *                          request.
794 *          VM_EVICT_FIRST : indicates to the virtual memory manager that when memory is low, this page should be
795 *                          evicted before other pages
796 * @param file If VM_FILE is set, then the page is backed by this file, starting at the position specified by pos.
797 *             If VM_FILE is clear, then this value must be NULL.
798 * @param pos  If VM_FILE is set, then this is the offset into the file where the page is mapped to. If VM_FILE is clear,
799 *             then this value must be 0.
800 *
801 * @maxirq1 IRQ1_SCHEDULER
802 */
803 static size_t MapVirtEx(struct vas* vas, size_t physical, size_t virtual, size_t pages, int flags, struct open_file* file, off_t pos) {
804     MAX_IRQ1_IRQ1_SCHEDULER;
805
806     /*
807      * We only specify a physical page when we need to map hardware directly (i.e. it's not
808      * part of the available RAM the physical memory manager can give).
809      */
810     if (physical != 0 && (flags & (VM_MAP_HARDWARE | VM_RELOCATABLE)) == 0) {
811         return 0;
812     }
813
814     if ((flags & VM_MAP_HARDWARE) && (flags & VM_FILE)) {
815         return 0;
816     }
817
818     if ((flags & VM_MAP_HARDWARE) && (flags & VM_LOCK) == 0) {
819         return 0;
820     }
821
822     if ((flags & VM_FILE) && file == NULL) {
823         return 0;
824     }
825
826     if ((flags & VM_FILE) == 0 && (file != NULL || pos != 0)) {
827         return 0;
828     }
829
830     if ((flags & VM_RELOCATABLE) && (flags & VM_FILE) == 0) {
831         return 0;
832     }
833
834     if ((flags & VM_RELOCATABLE) && physical == 0) {
835         return 0;
836     }
837
838     if ((flags & VM_FILE) && (flags & VM_LOCK)) {
839         return 0;
840     }
841
842     /*
843      * Get a virtual memory range that is not currently in use.
844      */
845     if (virtual == 0) {
846         virtual = AllocVirtRange(vas, pages, flags & VM_LOCAL);
847     } else {
848         // TODO: need to lock here to make the israngeinuse and allocvirtrange to be atomic
849         if (IsRangeInUse(vas, virtual, pages)) {
850             if (flags & VM_FIXED_VIRT) {
851                 return 0;
852             }
853         }
854         virtual = AllocVirtRange(vas, pages, flags & VM_LOCAL);
855     }
856
857 }
858
859 /*
860 * No point doing the multi-page mapping with only 2 pages, as the splitting cost is probably
861 * going to be greater than actually just adding 2 pages in the first place.
862 *
863 * May want to increase this value further in the future (e.g. maybe to 4 or 8)?
864 */
865 bool multi_page_mapping = ((flags & VM_LOCK) == 0) || ((flags & VM_MAP_HARDWARE) != 0) && pages >= 3;
866
867 for (size_t i = 0; i < (multi_page_mapping ? 1 : pages); ++i) {
868     if (flags & VM_FILE) {
869         ReferenceOpenFile(file);
870     }
871     AddMapping(
872         vas,
873         (flags & VM_RELOCATABLE) ? physical : (physical == 0 ? 0 : (physical + i * ARCH_PAGE_SIZE)),
874         virtual + i * ARCH_PAGE_SIZE,
875         flags,
876         file,
877         pos + i * ARCH_PAGE_SIZE,
878         multi_page_mapping ? pages : 1
879     );
880 }
881
882 if (vas == GetVas()) {
883     ArchFlushTlb(vas);
884 }
885
886 return virtual;
887 }
888
889 /**
890 * Creates a virtual memory mapping in the current virtual address space.
891 */
892 size_t MapVirt(size_t physical, size_t virtual, size_t bytes, int flags, struct open_file* file, off_t pos) {
893     return MapVirtEx(GetVas(), physical, virtual, BytesToPages(bytes), flags, file, pos);
894 }
895
896 static struct vas_entry* GetVirtEntry(struct vas* vas, size_t virtual) {
897     struct vas_entry dummy = { .num_pages = 1, .virtual = virtual & ~(ARCH_PAGE_SIZE - 1) };
898
899     assert(IsSpinlockHeld(&vas->lock));
900
901     struct vas_entry* res = (struct vas_entry*) AvlTreeGet(vas->mappings, (void*) &dummy);

```

```

902     if (res == NULL) {
903         AcquireSpinlockIrql(&GetCpu()->global_mappings_lock);
904         res = (struct vas_entry*) AvlTreeGet(GetCpu()->global_vas_mappings, (void*) &dummy);
905         ReleaseSpinlockIrql(&GetCpu()->global_mappings_lock);
906     }
907     return res;
908 }
909
910 size_t GetPhysFromVirt(size_t virtual) {
911     struct vas* vas = GetVas();
912     AcquireSpinlockIrql(&vas->lock);
913     struct vas_entry* entry = GetVirtEntry(vas, virtual);
914     size_t result = entry->physical;
915
916     /*
917     * Handle mappings of more than 1 page at a time by adding the extra offset
918     * from the start of the mapping.
919     */
920     size_t target_page = virtual / ARCH_PAGE_SIZE;
921     size_t entry_page = entry->virtual / ARCH_PAGE_SIZE;
922     if (entry_page < target_page) {
923         result += (target_page - entry_page) * ARCH_PAGE_SIZE;
924     }
925
926     ReleaseSpinlockIrql(&vas->lock);
927     return result;
928 }
929
930 static size_t SplitLargePageEntryIntoMultiple(struct vas* vas, size_t virtual, struct vas_entry* entry, int num_to_leave) {
931     if (entry->num_pages == 1) {
932         return 1;
933     }
934
935     if (entry->ref_count != 1) {
936         LogDeveloperWarning("Splitting multi-mapping with ref_count != 1, this hasn't been tested!\n");
937     }
938
939     /*
940     * Although it can't be allocated, it can be in RAM (e.g. for VM_MAP_HARDWARE).
941     */
942     assert(!entry->allocated);
943     assert(!entry->swapfile);
944
945     size_t entry_page = entry->virtual / ARCH_PAGE_SIZE;
946     size_t target_page = virtual / ARCH_PAGE_SIZE;
947
948     /*
949     * Split off anything before this page.
950     */
951     if (entry_page < target_page) {
952         size_t num_beforehand = target_page - entry_page;
953
954         struct vas_entry* pre_entry = AllocHeap(sizeof(struct vas_entry));
955         *pre_entry = *entry;
956
957         pre_entry->num_pages = num_beforehand;
958         entry->num_pages -= num_beforehand;
959         entry->virtual += num_beforehand * ARCH_PAGE_SIZE;
960
961         /*
962         * For multi-mapping for VM_MAP_HARDWARE
963         */
964         if (entry->physical != 0) {
965             entry->physical += num_beforehand * ARCH_PAGE_SIZE;
966         }
967
968         if (entry->file) {
969             entry->file_offset += num_beforehand * ARCH_PAGE_SIZE;
970         }
971
972         InsertIntoAvl(vas, pre_entry);
973     }
974
975     /*
976     * There's now no pages beforehand. Now we need to check if there are any other pages
977     * after this.
978     */
979     if (entry->num_pages > num_to_leave) {
980         struct vas_entry* post_entry = AllocHeap(sizeof(struct vas_entry));
981         *post_entry = *entry;
982
983         post_entry->num_pages -= num_to_leave;
984         entry->num_pages = num_to_leave;
985
986         post_entry->virtual += ARCH_PAGE_SIZE * num_to_leave;
987
988         /*
989         * For multi-mapping for VM_MAP_HARDWARE
990         */
991         if (entry->physical != 0) {
992             post_entry->physical += ARCH_PAGE_SIZE * num_to_leave;
993         }
994
995         if (entry->file) {
996             post_entry->file_offset += ARCH_PAGE_SIZE * num_to_leave;
997         }
998
999         InsertIntoAvl(vas, post_entry);
1000     }
1001
1002     return entry->num_pages;
1003 }
1004
1005 static void BringIntoMemoryFromCow(struct vas_entry* entry) {
1006     /*
1007     * If someone deallocates a COW page in another process to get the ref
1008     * count back to 1 already, then we just have the page to ourselves again.
1009     */
1010     if (entry->ref_count == 1) {
1011         entry->cow = false;
1012         ArchUpdateMapping(GetVas(), entry);
1013         ArchFlushTlb(GetVas());
1014         return;
1015     }
1016
1017     uint8_t page_data[ARCH_PAGE_SIZE];
1018     inline_memcpy(page_data, (void*) entry->virtual, ARCH_PAGE_SIZE);
1019
1020     entry->ref_count--;
1021
1022     if (entry->ref_count == 1) {
1023         entry->cow = false;
1024     }
1025
1026     struct vas_entry* new_entry = AllocHeap(sizeof(struct vas_entry));
1027     *new_entry = *entry;
1028     new_entry->ref_count = 1;
1029     new_entry->physical = AllocPhys();
1030     new_entry->allocated = true;
1031     DeleteFromAvl(GetVas(), entry);

```

```

1032     FreeHeap(entry);
1033     ArchUpdateMapping(GetVas(), entry);
1034     ArchFlushTlb(GetVas());
1035     inline_memcpy((void*) entry->virtual, page_data, ARCH_PAGE_SIZE);
1036 }
1037
1038 static void BringIntoMemoryFromFile(struct vas_entry* entry, size_t faulting_virt) {
1039     // TODO: need to test that you're allowed to read past the end of the file (even into other pages)
1040     // if the size mapped allows it, and just get zeros
1041
1042     SplitLargePageEntryIntoMultiple(GetVas(), faulting_virt, entry, 1);
1043     entry->load_in_progress = true;
1044     ArchUpdateMapping(GetVas(), entry);
1045     ArchFlushTlb(GetVas());
1046     DeferDiskRead(entry->virtual, entry->file_node, entry->file_offset, false);
1047 }
1048
1049 static void BringIntoMemoryFromSwapfile(struct vas_entry* entry) {
1050     assert(!entry->file);
1051
1052     uint64_t offset = entry->swapfile_offset;
1053     entry->load_in_progress = true;
1054     ArchUpdateMapping(GetVas(), entry);
1055     ArchFlushTlb(GetVas());
1056     DeferDiskRead(entry->virtual, GetSwapfile(), offset, true);
1057 }
1058
1059 static void BringInBlankPage(struct vas* vas, struct vas_entry* entry, size_t faulting_virt, int fault_type) {
1060     if ((fault_type & VM_READ) && !entry->read) {
1061         UnhandledFault();
1062     }
1063     if ((fault_type & VM_WRITE) && !entry->write) {
1064         UnhandledFault();
1065     }
1066     if ((fault_type & VM_EXEC) && !entry->exec) {
1067         UnhandledFault();
1068     }
1069
1070     SplitLargePageEntryIntoMultiple(vas, faulting_virt, entry, 1);
1071     assert(entry->num_pages == 1);
1072
1073     entry->physical = AllocPhys();
1074     entry->allocated = true;
1075     entry->in_ram = true;
1076     entry->allow_temp_write = true;
1077     assert(!entry->swapfile);
1078     ArchUpdateMapping(vas, entry);
1079     ArchFlushTlb(vas);
1080
1081     inline_memset((void*) entry->virtual, 0, ARCH_PAGE_SIZE);
1082     entry->allow_temp_write = false;
1083     ArchUpdateMapping(vas, entry);
1084     ArchFlushTlb(vas);
1085 }
1086
1087 static int BringIntoMemory(struct vas* vas, struct vas_entry* entry, bool allow_cow, size_t faulting_virt, int fault_type) {
1088     assert(IsSpinlockHeld(&vas->lock));
1089
1090     if (entry->cow && allow_cow) {
1091         assert(entry->num_pages == 1);
1092         BringIntoMemoryFromCow(entry);
1093         return 0;
1094     }
1095
1096     if (entry->file && !entry->in_ram) {
1097         BringIntoMemoryFromFile(entry, faulting_virt);
1098         return 0;
1099     }
1100
1101     if (entry->swapfile) {
1102         assert(entry->num_pages == 1);
1103         BringIntoMemoryFromSwapfile(entry);
1104         return 0;
1105     }
1106
1107     if (!entry->in_ram) {
1108         BringInBlankPage(vas, entry, faulting_virt, fault_type);
1109         return 0;
1110     }
1111
1112     return EINVAL;
1113 }
1114
1115 bool LockVirtEx(struct vas* vas, size_t virtual) {
1116     struct vas_entry* entry = GetVirtEntry(vas, virtual);
1117
1118     if (!entry->in_ram) {
1119         SplitLargePageEntryIntoMultiple(vas, virtual, entry, 1);
1120         int res = BringIntoMemory(vas, entry, true, virtual, 0);
1121         if (res != 0) {
1122             Panic(PANIC_CANNOT_LOCK_MEMORY);
1123         }
1124         assert(entry->in_ram);
1125     }
1126
1127     bool old_lock = entry->lock;
1128     entry->lock = true;
1129     return old_lock;
1130 }
1131
1132 void UnlockVirtEx(struct vas* vas, size_t virtual) {
1133     struct vas_entry* entry = GetVirtEntry(vas, virtual);
1134     SplitLargePageEntryIntoMultiple(vas, virtual, entry, 1);
1135     entry->lock = false;
1136 }
1137
1138 bool LockVirt(size_t virtual) {
1139     struct vas* vas = GetVas();
1140     AcquireSpinlockIrql(&vas->lock);
1141     bool res = LockVirtEx(vas, virtual);
1142     ReleaseSpinlockIrql(&vas->lock);
1143     return res;
1144 }
1145
1146 void UnlockVirt(size_t virtual) {
1147     struct vas* vas = GetVas();
1148     AcquireSpinlockIrql(&vas->lock);
1149     UnlockVirtEx(vas, virtual);
1150     ReleaseSpinlockIrql(&vas->lock);
1151 }
1152
1153 void SetVirtPermissions(size_t virtual, int set, int clear) {
1154     /*
1155     * Only allow these flags to be set / cleared.
1156     */
1157     if ((set | clear) & ~(VM_READ | VM_WRITE | VM_EXEC | VM_USER)) {
1158         assert(false);
1159         return;
1160     }
1161 }

```

```

1162 struct vas* vas = GetVas();
1163 AcquireSpinlockIrql(&vas->lock);
1164
1165 struct vas_entry* entry = GetVirtEntry(vas, virtual);
1166 if (entry == NULL)
1167     PanicEx(PANIC_ASSERTION_FAILURE, "[SetVirtPermissions] got null back for virt_entry");
1168 }
1169
1170 SplitLargePageEntryIntoMultiple(vas, virtual, entry, 1);
1171
1172 entry->read = (set & VM_READ ? true : (clear & VM_READ ? false : entry->read);
1173 entry->write = (set & VM_WRITE ? true : (clear & VM_WRITE ? false : entry->write);
1174 entry->exec = (set & VM_EXEC ? true : (clear & VM_EXEC ? false : entry->exec);
1175 entry->user = (set & VM_USER ? true : (clear & VM_USER ? false : entry->user);
1176
1177 ArchUpdateMapping(vas, entry);
1178 ArchFlushTlb(vas);
1179
1180 ReleaseSpinlockIrql(&vas->lock);
1181 }
1182
1183 int GetVirtPermissions(size_t virtual) {
1184     struct vas* vas = GetVas();
1185     AcquireSpinlockIrql(&vas->lock);
1186     struct vas_entry* entry_ptr = GetVirtEntry(GetVas(), virtual);
1187     if (entry_ptr == NULL)
1188         ReleaseSpinlockIrql(&vas->lock);
1189     return 0;
1190 }
1191 struct vas_entry entry = *entry_ptr;
1192 ReleaseSpinlockIrql(&vas->lock);
1193
1194 int permissions = 0;
1195 if (entry.read) permissions |= VM_READ;
1196 if (entry.write) permissions |= VM_WRITE;
1197 if (entry.exec) permissions |= VM_EXEC;
1198 if (entry.lock) permissions |= VM_LOCK;
1199 if (entry.file) permissions |= VM_FILE;
1200 if (entry.user) permissions |= VM_USER;
1201 if (entry.global) permissions |= VM_LOCAL;
1202 if (entry.relocatable) permissions |= VM_RELOCATABLE;
1203
1204 return permissions;
1205 }
1206
1207 int UnmapVirtEx(struct vas* vas, size_t virtual, size_t pages) {
1208     bool needs_tlb_flush = false;
1209
1210     for (size_t i = 0; i < pages; ++i) {
1211         struct vas_entry* entry = GetVirtEntry(vas, virtual + i * ARCH_PAGE_SIZE);
1212         if (entry == NULL) {
1213             return EINVAL;
1214         }
1215
1216         SplitLargePageEntryIntoMultiple(vas, virtual, entry, 1); // TODO: multi-pages
1217
1218         assert(entry->ref_count > 0);
1219         entry->ref_count--;
1220
1221         if (entry->ref_count == 0) {
1222             if (entry->file && entry->write && entry->in_ram) {
1223                 DeferDiskWrite(entry->virtual, entry->file_node, entry->file_offset);
1224                 // TODO: after that DeferDiskWrite, we need to defer a DereferenceOpenFile(entry->file_node)
1225             }
1226             if (entry->in_ram) {
1227                 ArchUnmap(vas, entry);
1228                 needs_tlb_flush = true;
1229             }
1230             if (entry->swapfile) {
1231                 assert(!entry->allocated);
1232                 DeallocateswapfileIndex(entry->physical / ARCH_PAGE_SIZE);
1233             }
1234             if (entry->allocated) {
1235                 assert(!entry->swapfile); // can't be on swap, as putting on swap clears allocated bit
1236                 DeallocPhys(entry->physical);
1237             }
1238
1239             DeleteFromAvl(vas, entry);
1240             FreeHeap(entry);
1241             FreeVirtRange(vas, virtual + i * ARCH_PAGE_SIZE, entry->num_pages);
1242         }
1243     }
1244
1245     if (needs_tlb_flush) {
1246         ArchFlushTlb(vas);
1247     }
1248
1249     return 0;
1250 }
1251
1252 int UnmapVirt(size_t virtual, size_t bytes) {
1253     struct vas* vas = GetVas();
1254     AcquireSpinlockIrql(&vas->lock);
1255     int res = UnmapVirtEx(vas, virtual, BytesToPages(bytes));
1256     ReleaseSpinlockIrql(&vas->lock);
1257     return res;
1258 }
1259
1260 static void CopyVasRecursive(struct avl_node* node, struct vas* new_vas) {
1261     if (node == NULL) {
1262         return;
1263     }
1264
1265     CopyVasRecursive(AvlTreeGetLeft(node), new_vas);
1266     CopyVasRecursive(AvlTreeGetRight(node), new_vas);
1267
1268     struct vas_entry* entry = AvlTreeGetData(node);
1269
1270     if (entry->lock) {
1271         /*
1272          * Got to add the new entry right now. We know it must be in memory as it
1273          * is locked.
1274          */
1275         assert(entry->in_ram);
1276
1277         if (entry->allocated) {
1278             /*
1279              * Copy the physical page. We do this by copying the data into a buffer,
1280              * putting a new physical page in the existing VAS and then copying the
1281              * data there. Then the original physical page that was there is free to use
1282              * as the copy.
1283              */
1284             uint8_t page_data[ARCH_PAGE_SIZE]; // TODO: MapVirt this ?
1285             inline_memcpy(page_data, (void*) entry->virtual, ARCH_PAGE_SIZE);
1286             size_t new_physical = entry->physical;
1287             entry->physical = AllocPhys();
1288             ArchUpdateMapping(GetVas(), entry);
1289             ArchFlushTlb(GetVas());
1290             inline_memcpy((void*) entry->virtual, page_data, ARCH_PAGE_SIZE);
1291         }
1292     }

```

```

1292     struct vas_entry* new_entry = AllocHeap(sizeof(struct vas_entry));
1293     *new_entry = *entry;
1294     new_entry->ref_count = 1;
1295     new_entry->physical = new_physical;
1296     new_entry->allocated = true;
1297     AvlTreeInsert(new_vas->mappings, entry); // don't need to insert global - we're copying so it's already in global
1298     ArchAddMapping(new_vas, entry);
1299
1300 } else {
1301     LogWriteSerial("fork() on a hardware-mapped page is not implemented yet");
1302     PanicEx(PANIC_NOT_IMPLEMENTED, "CopyVasRecursive");
1303 }
1304
1305 } else {
1306     /*
1307     * If it's on swap, it's okay to still mark it as COW, as when we reload we will
1308     * try to do the 'copy'-on-write, and then we will reload from swap, and it will
1309     * then reload and then be copied. Alternatively, if it is read, then it gets brought
1310     * back into memory, but as a COW page still.
1311     */
1312     * BSS memory works fine like this too (but will incur another fault when it is used).
1313
1314     * At this stage (where shared memory doesn't exist yet), file mapped pages will also
1315     * be COWed. This means there will be two copies of the file in memory should they write
1316     * to it. The final process to release memory will ultimately 'win' and have its changes
1317     * perserved to disk (the others will get overwritten).
1318     */
1319     entry->cow = true;
1320     entry->ref_count++;
1321
1322     // again, no need to add to global - it's already there!
1323     AvlTreeInsert(new_vas->mappings, entry);
1324
1325     ArchUpdateMapping(GetVas(), entry);
1326     ArchAddMapping(new_vas, entry);
1327 }
1328 }
1329
1330 struct vas* CopyVas(void) {
1331     struct vas* vas = GetVas();
1332     struct vas* new_vas = CreateVas();
1333
1334     AcquireSpinlockIrql(&vas->lock);
1335     // no need to change global - it's already there!
1336     CopyVasRecursive(AvlTreeGetRootNode(vas->mappings), new_vas);
1337     ArchFlushTlb(vas);
1338     ReleaseSpinlockIrql(&vas->lock);
1339
1340     return new_vas;
1341 }
1342
1343 struct vas* GetVas(void) {
1344     // TODO: cpu probably needs to have a lock object in it called current_vas_lock, which needs to be held whenever
1345     // someone reads or writes to current_vas;
1346     return GetCpu()->current_vas;
1347 }
1348
1349 void SetVas(struct vas* vas) {
1350     GetCpu()->current_vas = vas;
1351     ArchSetVas(vas);
1352 }
1353
1354 struct vas* GetKernelVas(void) {
1355     return kernel_vas;
1356 }
1357
1358 void InitVirt(void) {
1359     // TODO: cpu probably needs to have a lock object in it called current_vas_lock, which needs to be held whenever
1360     // someone reads or writes to current_vas;
1361
1362     assert(!virt_initialised);
1363     GetCpu()->global_vas_mappings = AvlTreeCreate();
1364     AvlTreeSetComparator(GetCpu()->global_vas_mappings, VirtAvlComparator);
1365     ArchInitVirt();
1366
1367     kernel_vas = GetVas();
1368     virt_initialised = true;
1369 }
1370
1371 /**
1372  * Handles a page fault. Only to be called by the low-level, platform specific interrupt handler when a page
1373  * fault occurs. It will attempt to resolve any fault (e.g. handling copy-on-write, swapfile, file-backed, etc.).
1374  */
1375 * @param faulting_virt The virtual address that was accessed that caused the page fault
1376 * @param fault_type The reason why a page fault occurred. Is a bitfield of VM_WRITE, VM_READ, VM_USER and VM_EXEC.
1377 * VM_READ should be set if a non-present page was accessed. VM_USER should be set for permission
1378 * faults, and VM_WRITE should be set if the operation was caused by a write (as opposed to a read).
1379 * VM_EXEC should be set if execution tried to occur in a non-executable page.
1380 */
1381 * @maxirql IRQL_PAGE_FAULT
1382 */
1383 int handling_page_fault = 0;
1384
1385 void HandleVirtFault(size_t faulting_virt, int fault_type) {
1386     if (GetIrql() >= IRQL_SCHEDULER) {
1387         PanicEx(PANIC_INVALID_IRQL, "page fault while IRQL >= IRQL_SCHEDULER. is some clown holding a spinlock while "
1388             "executing pageable code? or calling AllocHeapEx wrong with a lock held?");
1389     }
1390
1391     struct vas* vas = GetVas();
1392     AcquireSpinlockIrql(&vas->lock);
1393     ++handling_page_fault;
1394
1395     struct vas_entry* entry = GetVirtEntry(vas, faulting_virt);
1396
1397     if (entry == NULL) {
1398         UnhandledFault();
1399     }
1400
1401     if (entry->load_in_progress) {
1402         --handling_page_fault;
1403         ReleaseSpinlockIrql(&vas->lock);
1404         Schedule();
1405         return;
1406     }
1407
1408     /*
1409     * Sanity check that our flags are configured correctly.
1410     */
1411     assert(!(entry->in_ram && entry->swapfile));
1412     assert(!(entry->file && entry->swapfile));
1413     assert(!(entry->in_ram && entry->lock));
1414     assert(!(entry->cow && entry->lock));
1415
1416     // TODO: check for access violations (e.g. user using a supervisor page)
1417     // (read / write is not necessarily a problem, e.g. COW)
1418
1419     int result = BringIntoMemory(vas, entry, fault_type & VM_WRITE, faulting_virt, fault_type);
1420     if (result != 0) {

```

```

1422     UnhandledFault();
1423 }
1424
1425 --handling_page_fault;
1426 ReleaseSpinlockIrql(&vas->lock);
1427 }
1428
1429 /**
1430 * Determines whether or not virtual memory has been initialised yet. This can be used to determine if it
1431 * is possible to call any virtual memory functions (e.g. in the physical memory and heap allocators).
1432 *
1433 * @return True if virtual memory is available, false otherwise.
1434 *
1435 * @maxirql IRQL_HIGH
1436 */
1437 bool IsVirtInitialised(void) {
1438     return virt_initialised;
1439 }
1440
1441 size_t BytesToPages(size_t bytes) {
1442     return (bytes + ARCH_PAGE_SIZE - 1) / ARCH_PAGE_SIZE;
1443 }
1444
1445 void DestroyVas(struct vas* vas) {
1446     /*
1447     * TODO: implement this
1448     */
1449     (void) vas;
1450
1451     if (vas == GetVas()) {
1452         Panic(PANIC_VAS_TRIED_TO_SELF_DESTRUCT);
1453     }
1454
1455     // TODO: may need to add reference counting later on (depending on what we need),
1456     // just decrement here, and only delete if got to 0.
1457 }

```