

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 (KRNAPI.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  * - VnodeOpWait, select/poll syscalls
46  */
47
48 void DummyAppThread(void*) {
49     PutsConsole("drv0:> ");
50
51     struct open_file* con;
52     OpenFile("con:", O_RDONLY, 0, &con);
53
54     while (true) {
55         char bf[302];
56         inline memset(bf, 0, 302);
57         struct transfer tr = CreateKernelTransfer(bf, 301, 0, TRANSFER_READ);
58         ReadFile(con, &tr);
59         PutsConsole("Command not found: ");
60         PutsConsole(bf);
61         PutsConsole("\n");
62
63         if (bf[0] == 'u' || bf[0] == 'U') {
64             CreateUserModeProcess(NULL, "sys:/init.exe");
65         }
66         else if (bf[0] == 'p' || bf[0] == 'P') {
67             Panic(PANIC_MANUALLY_INITIATED);
68         }
69         else if (bf[0] == 'e' || bf[0] == 'E') {
70             MapVirt(0, 0, 8 * 4096, VM_LOCK | VM_READ, NULL, 0);
71         }
72
73         PutsConsole("drv0:> ");
74     }
75 }
76
77 void InitUserspace(void) {
78     size_t free = GetFreePhysKilobytes();
79     size_t total = GetTotalPhysKilobytes();
80     DbgScreenPrintf("NOS Kernel\nCopyright Alex Boxall 2022-2023\n\n%ld / %ld KB used (%ld%% free)\n\n", total - free, total, 100 * (free) / total);
81     CreateThread(DummyAppThread, NULL, GetVas(), "dummy app");
82 }
83
84 void InitThread(void*) {
85     ReinitHeap();
86     InitRandomDevice();
87     InitNullDevice();
88     InitConsole();
89     InitProcess();
90     InitDiskCaches();
91
92     InitFilesystemTable();
93     ArchInitDev(false);
94
95     struct open_file* sys_folder;
96     int res = OpenFile("drv0:/System", O_RDONLY, 0, &sys_folder);
97     if (res != 0) {
98         PanicEx(PANIC_NO_FILESYSTEM, "sys A");
99     }
100     res = AddVfsMount(sys_folder->node, "sys");
101     if (res != 0) {
102         PanicEx(PANIC_NO_FILESYSTEM, "sys B");
103     }
104
105     struct open_file* swapfile;
106     res = OpenFile("raw-hd0:/part1", O_RDWR, 0, &swapfile);
107     if (res != 0) {
108         PanicEx(PANIC_NO_FILESYSTEM, "swapfile A");
109     }
110     res = AddVfsMount(swapfile->node, "swap");
111     if (res != 0) {
112         PanicEx(PANIC_NO_FILESYSTEM, "swapfile B");
113     }
114
115     InitSwapfile();
116     InitSymbolTable();
117     ArchInitDev(true);
118     InitProgramLoader();
119     InitUserspace();
120
121     MarkTfwStartPoint(TFW_SP_ALL_CLEAR);
122
123     while (true) {
124         /*
125          * We crash in strange and rare conditions if this thread's stack gets removed, so we will
126          * ensure we don't terminate it.
127          */
128         SleepMilli(100000);
129     }
130 }
```

```

129     }
130 }
131
132 #include <machine/portio.h>
133 static void InitSerialDebugging(void) {
134     const int PORT = 0x3F8;
135     outb(PORT + 1, 0x00); // Disable all interrupts
136     outb(PORT + 3, 0x80); // Enable DLAB (set baud rate divisor)
137     outb(PORT + 0, 0x03); // Set divisor to 3 (lo byte) 38400 baud
138     outb(PORT + 1, 0x00); // (hi byte)
139     outb(PORT + 3, 0x03); // 8 bits, no parity, one stop bit
140     outb(PORT + 2, 0xC7); // Enable FIFO, clear them, with 14-byte threshold
141     outb(PORT + 4, 0x0B); // IRQs enabled, RTS/DSR set
142     outb(PORT + 4, 0x1E); // Set in loopback mode, test the serial chip
143     outb(PORT + 0, 0xAE); // Test serial chip (send byte 0xAE and check if serial returns same byte)
144
145     // Check if serial is faulty (i.e: not same byte as sent)
146     if (inb(PORT + 0) != 0xAE) {
147         return;
148     }
149
150     // If serial is not faulty set it in normal operation mode
151     // (not-loopback with IRQs enabled and OUT#1 and OUT#2 bits enabled)
152     outb(PORT + 4, 0x0F);
153 }
154
155 void KernelMain(void) {
156     InitSerialDebugging();
157
158     LogWriteSerial("KernelMain: kernel is initialising...\n");
159
160     /*
161      * Allows us to call GetCpu(), which allows IRQL code to work. Anything which uses
162      * IRQL (i.e. the whole system) relies on this, so this must be done first.
163      */
164     InitCpuTable();
165     assert(GetIrql() == IRQL_STANDARD);
166
167     /*
168      * Initialise the testing framework if we're in debug mode.
169      */
170     InitTfw();
171     MarkTfwStartPoint(TFW_SP_INITIAL);
172
173     InitPhys();
174     MarkTfwStartPoint(TFW_SP_AFTER_PHYS);
175
176     /*
177      * Allows deferments of functions to actually happen. IRQL is still usable beforehand though.
178      */
179     InitIrql();
180     InitVfs();
181     InitTimer();
182     InitScheduler();
183     InitDiskUtil();
184
185     InitHeap();
186     MarkTfwStartPoint(TFW_SP_AFTER_HEAP);
187
188     InitBootstrapCpu();
189     MarkTfwStartPoint(TFW_SP_AFTER_BOOTSTRAP_CPU);
190
191     InitVirt();
192     MarkTfwStartPoint(TFW_SP_AFTER_VIRT);
193
194     ReinitPhys();
195     MarkTfwStartPoint(TFW_SP_AFTER_PHYS_REINIT);
196
197     InitOtherCpu();
198     MarkTfwStartPoint(TFW_SP_AFTER_ALL_CPU);
199
200     CreateThreadEx(InitThread, NULL, GetVas(), "init", NULL, SCHEDULE_POLICY_FIXED, FIXED_PRIORITY_KERNEL_NORMAL, 0);
201     StartMultitasking();
202 }
203

```

File: `/adt/threadlist.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     else if (bf == 2) {
75         if (AvlGetBalance(tree->left) == -1) {
76             tree->left = AvlRotateLeft(tree->left);
77         }
78         return AvlRotateRight(tree);
79     }
80     else {
81         return tree;
82     }
83 }
84
85 static struct avl_node* AvlInsert(struct avl_node* tree, void* data, avl_comparator comparator) {
86     struct avl_node* new_tree;
87
88     assert(comparator != NULL);
89     assert(tree != NULL);
90
91     if (comparator(data, tree->data) < 0) {
92         struct avl_node* left_tree;
93         if (tree->left == NULL) {
94             left_tree = AvlCreateNode(data, NULL, NULL);
95         } else {
96             left_tree = AvlInsert(tree->left, data, comparator);
97         }
98
99         new_tree = AvlCreateNode(tree->data, left_tree, tree->right);
100     } else {
101         struct avl_node* right_tree;
102         if (tree->right == NULL) {
103             right_tree = AvlCreateNode(data, NULL, NULL);
104         } else {
105             right_tree = AvlInsert(tree->right, data, comparator);
106         }
107
108         new_tree = AvlCreateNode(tree->data, tree->left, right_tree);
109     }
110
111     FreeHeap(tree);
112
113     return AvlBalance(new_tree);
114 }
115
116 static struct avl_node* AvlDelete(struct avl_node* tree, void* data, avl_comparator comparator) {
117     if (tree == NULL) {
118         return NULL;
119     }
120
121     struct avl_node* to_free = NULL;
122
123     if (comparator(data, tree->data) < 0) {
124         tree->left = AvlDelete(tree->left, data, comparator);
125     }
126     else if (comparator(data, tree->data) > 0) {
127         tree->right = AvlDelete(tree->right, data, comparator);
128     }
129     else if (tree->left == NULL) {
130         to_free = tree;
131         tree = tree->right;
132     }
133     else if (tree->right == NULL) {
134         to_free = tree;
135         tree = tree->left;
136     }
137     else {
138         struct avl_node* node = tree->right;
139         while (node->left != NULL) {
140             node = node->left;
141         }
142
143         tree->data = node->data;
144         tree->right = AvlDelete(tree->right, node->data, comparator);
145     }
146
147     /*
148     * If NULL is passed in to FreeHeap, nothing happens (which is what we want).
149     */
150

```

```

151     FreeHeap(to_free);
152
153     return AvlBalance(tree);
154 }
155
156 /**
157  * Given an object, find it in the AVL tree and return it. This is useful if the comparator only compares
158  * part of the object, and so the entire object can be retrieved by searching for only part of it.
159  */
160 static void* AvlGet(struct avl_node* tree, void* data, avl_comparator comparator) {
161     if (tree == NULL) {
162         return NULL;
163     }
164
165     if (comparator(tree->data, data) == 0) {
166         /*
167          * Must return 'tree->data', (and not 'data'), as tree->data != data if there is a custom comparator.
168          */
169         return tree->data;
170     }
171
172     void* left = AvlGet(tree->left, data, comparator);
173     if (left != NULL) {
174         return left;
175     }
176     return AvlGet(tree->right, data, comparator);
177 }
178
179 static void AvlPrint(struct avl_node* tree, void(*printer)(void*)) {
180     if (tree == NULL) {
181         return;
182     }
183     AvlPrint(tree->left, printer);
184     if (printer == NULL) {
185         LogWriteSerial("[[0x%X]], \n", tree->data);
186     } else {
187         printer(tree->data);
188     }
189     AvlPrint(tree->right, printer);
190 }
191
192 static bool AvlContains(struct avl_node* tree, void* data, avl_comparator comparator) {
193     if (tree == NULL) {
194         return false;
195     }
196
197     if (comparator(tree->data, data) == 0) {
198         return true;
199     }
200
201     return AvlContains(tree->left, data, comparator) || AvlContains(tree->right, data, comparator);
202 }
203
204 static void AvlDestroy(struct avl_node* tree, avl_deletion_handler handler) {
205     if (tree == NULL) {
206         return;
207     }
208
209     AvlDestroy(tree->left, handler);
210     AvlDestroy(tree->right, handler);
211     if (handler != NULL) {
212         handler(tree->data);
213     }
214     FreeHeap(tree);
215 }
216
217 static int AvlDefaultComparator(void* a, void* b) {
218     if (a == b) return 0;
219     return (a < b) ? -1 : 1;
220 }
221
222 struct avl_tree* AvlTreeCreate(void) {
223     struct avl_tree* tree = AllocHeap(sizeof(struct avl_tree));
224     tree->size = 0;
225     tree->root = NULL;
226     tree->deletion_handler = NULL;
227     tree->equality_handler = AvlDefaultComparator;
228     return tree;
229 }
230
231 avl_deletion_handler AvlTreeSetDeletionHandler(struct avl_tree* tree, avl_deletion_handler handler) {
232     avl_deletion_handler ret = tree->deletion_handler;
233     tree->deletion_handler = handler;
234     return ret;
235 }
236
237 avl_comparator AvlTreeSetComparator(struct avl_tree* tree, avl_comparator handler) {
238     avl_comparator ret = tree->equality_handler;
239     tree->equality_handler = handler;
240     return ret;
241 }
242
243 void AvlTreeInsert(struct avl_tree* tree, void* data) {
244     if (tree->root == NULL) {
245         tree->root = AvlCreateNode(data, NULL, NULL);
246     } else {
247         tree->root = AvlInsert(tree->root, data, tree->equality_handler);
248     }
249     tree->size++;
250 }
251
252 void AvlTreeDelete(struct avl_tree* tree, void* data) {
253     tree->root = AvlDelete(tree->root, data, tree->equality_handler);
254     tree->size--;
255 }
256
257 bool AvlTreeContains(struct avl_tree* tree, void* data) {
258     return AvlContains(tree->root, data, tree->equality_handler);
259 }
260
261 void* AvlTreeGet(struct avl_tree* tree, void* data) {
262     return AvlGet(tree->root, data, tree->equality_handler);
263 }
264
265 int AvlTreeSize(struct avl_tree* tree) {
266     return tree->size;
267 }
268
269 void AvlTreeDestroy(struct avl_tree* tree) {
270     AvlDestroy(tree->root, tree->deletion_handler);
271     FreeHeap(tree);
272 }
273
274 struct avl_node* AvlTreeGetRootNode(struct avl_tree* tree) {
275     return tree->root;
276 }
277
278 struct avl_node* AvlTreeGetLeft(struct avl_node* node) {
279     if (node == NULL) {
280         return NULL;

```

```

281     }
282
283     return node->left;
284 }
285
286 struct avl_node* AvlTreeGetRight(struct avl_node* node) {
287     if (node == NULL) {
288         return NULL;
289     }
290
291     return node->right;
292 }
293
294 void* AvlTreeGetData(struct avl_node* node) {
295     if (node == NULL) {
296         return NULL;
297     }
298
299     return node->data;
300 }
301
302 void AvlTreePrint(struct avl_tree* tree, void(*printer)(void*)) {
303     AvlPrint(tree->root, printer);
304 }

```

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     if (right < queue->size) {
66         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])) {
67             extreme = right;
68         }
69     }
70     if (i != extreme) {
71         SwapElements(queue, i, extreme);
72         Heapify(queue, extreme);
73     }
74 }
75
76 void PriorityQueueInsert(struct priority_queue* queue, void* elem, uint64_t priority) {
77     if (queue->size == queue->capacity) {
78         // I think this can happen when the OS is running too slowly! (the deferred function buffer actually fills up
79         // and overflows). Testing on real H/W from 1996, debug prints took about a second per character!!
80         PanicEx(PANIC_PRIORITY_QUEUE, "insert called when full");
81     }
82
83     int i = queue->size++;
84     queue->array[i * queue->qwords_per_element] = priority;
85     inline_memcpy(queue->array + i * queue->qwords_per_element + 1, elem, queue->element_width);
86
87     if (queue->max) {
88         while (i != 0 && queue->array[(i - 1) / 2] * queue->qwords_per_element < queue->array[i * queue->qwords_per_element]) {
89             SwapElements(queue, (i - 1) / 2, i);
90             i = (i - 1) / 2;
91         }
92     } else {
93         while (i != 0 && queue->array[(i - 1) / 2] * queue->qwords_per_element > queue->array[i * queue->qwords_per_element]) {
94             SwapElements(queue, (i - 1) / 2, i);
95             i = (i - 1) / 2;
96         }
97     }
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 IRQL 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 <irq.h>
6
7  /**
8   * Creates a new open file from an opened vnode. An open file is used to link a vnode with corresponding data,
9   * about a particular instance of opening a file: such as a seek position, and ability to read or write; and is
10   * used to maintain file descriptor tables for the C userspace library.
11   *
12   * Open files maintain a reference counter that can be incremented and decremented with ReferenceOpenFile and
13   * DereferenceOpenFile. A newly created open file has a reference count of 1. When the count reaches 0, the memory
14   * is freed.
15   *
16   * @param node      The already-open vnode to wrap with additional data
17   * @param mode       The Unix-permissions passed to the OpenFile when the vnode was opened. Ignored by the kernel
18   *                  so far.
19   * @param flag       The flags that were passed to OpenFile when the vnode was opened. Should be a bitfield consisting of
20   *                  zero or more of: O_RDONLY, O_WRONLY, O_TRUNC, O_CREAT. See OpenFile for details. The storage of these
21   *                  flags in an open file should only be of interest to usermode programs - the flags here may be zero
22   *                  if the open file was created within the kernel, so the flags should be ignored within the kernel.
23   * @param can_read   Whether or not this open file is allowed to make read operations on the underlying vnode
24   * @param can_write  Whether or not this open file is allowed to make write operations on the underlying vnode
25   *
26   * @return A pointer to the newly created open file.
27   */
28 struct open_file* CreateOpenFile(struct vnode* node, int mode, int flags, bool can_read, bool can_write) {
29     MAX_IRQL(IRQL_SCHEDULER);
30
31     struct open_file* file = AllocHeap(sizeof(struct open_file));
32     file->reference_count = 1;
33     file->node = node;
34     file->can_read = can_read;
35     file->can_write = can_write;
36     file->initial_mode = mode;
37     file->flags = flags;
38     file->seek_position = 0;
39     InitSpinlock(&file->reference_count_lock, "open file", IRQL_SCHEDULER);
40
41     return file;
42 }
43
44 /**
45   * Increments the reference counter for an opened file. This should be called everytime a reference to
46   * the open file is kept, so that its memory can be managed correctly.
47   *
48   * @param file The open file to reference
49   */
50 void ReferenceOpenFile(struct open_file* file) {
51     MAX_IRQL(IRQL_SCHEDULER);
52     assert(file != NULL);
53
54     AcquireSpinlockIrql(&file->reference_count_lock);
55     file->reference_count++;
56     ReleaseSpinlockIrql(&file->reference_count_lock);
57 }
58
59 /**
60   * Decrements the reference counter for an opened file. Should be called whenever a reference to the open
61   * file is removed. If the reference counter reaches zero, the memory behind the open file will be freed.
62   * The underlying vnode within the open file is not dereferenced - this should be done prior to calling this
63   * function.
64   *
65   * @param file The open file to dereference
66   */
67 void DereferenceOpenFile(struct open_file* file) {
68     MAX_IRQL(IRQL_SCHEDULER);
69     assert(file != NULL);
70
71     AcquireSpinlockIrql(&file->reference_count_lock);
72
73     assert(file->reference_count > 0);
74     file->reference_count--;
75
76     if (file->reference_count == 0) {
77         /*
78          * Must release the lock before we delete it so we can put interrupts back on
79          */
80         ReleaseSpinlockIrql(&file->reference_count_lock);
81         FreeHeap(file);
82         return;
83     }
84
85     ReleaseSpinlockIrql(&file->reference_count_lock);
86 }

```

File: /vfs/diskutil.c

```

1  #include <diskutil.h>
2  #include <string.h>
3  #include <irq.h>
4  #include <assert.h>
5  #include <spinlock.h>
6  #include <vfs.h>
7  #include <errno.h>
8  #include <log.h>
9  #include <irq.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     } else if (num < 100) {
70         num_str[0] = (num / 10) + '0';
71         num_str[1] = (num % 10) + '0';
72     } else {
73         num_str[0] = (num / 100) + '0';
74         num_str[1] = ((num / 10) % 10) + '0';
75         num_str[2] = (num % 10) + '0';
76     }
77
78     strcat(str, num_str);
79     return 0;
80 }
81
82 /**
83  * Returns the name the next-mounted filesystem should receive (e.g. drv0,
84  * drv1, etc.) Each call to this function will return a different string. The
85  * caller is responsible for freeing the returned string.
86  *
87  * @return A caller-free string representing the drive name.
88  *
89  * @maxirql IRQL_SCHEDULER
90  */
91 char* GenerateNewMountedDiskName() {
92     MAX_IRQL(IRQL_SCHEDULER);
93
94     char name[16];
95     strcpy(name, "drv");
96
97     AcquireSpinlockIrql(&type_table_lock);
98     int disk_num = next_mounted_disk_num++;
99     ReleaseSpinlockIrql(&type_table_lock);
100
101     AppendNumberToString(name, disk_num);
102     return strdup(name);
103 }
104
105 /**
106  * Returns the name the next-mounted raw disk should receive, based on its type (e.g. raw-hd0, raw-hd1,
107  * raw-fd0). Each call to this function will return a different string. The caller is responsible for
108  * freeing the returned string.
109  *
110  * @param type The type of disk (one of DISKUTIL_TYPE_...)
111  * @return The caller-free string representing the drive name.
112  *
113  * @maxirql IRQL_SCHEDULER
114  */
115 char* GenerateNewRawDiskName(int type) {
116     MAX_IRQL(IRQL_SCHEDULER);
117
118     char name[16];
119     strcpy(name, "raw-");
120
121     if (type >= __DISKUTIL_NUM_TYPES || type < 0) {
122         type = DISKUTIL_TYPE_OTHER;
123     }
124
125     strcat(name, type_strings[type]);
126
127     AcquireSpinlockIrql(&type_table_lock);
128     int disk_num = type_table[type]++;
129     ReleaseSpinlockIrql(&type_table_lock);
130
131     AppendNumberToString(name, disk_num);
132     return strdup(name);
133 }
134
135 /**
136  * Generates and returns the name of a partition from its partition index within a drive (e.g. part0, part1)
137  * The caller is responsible for freeing the returned string.
138  */
139 static char* GetPartitionNameString(int index) {
140     char name[16];
141     strcpy(name, "part");
142     AppendNumberToString(name, index);
143     return strdup(name);
144 }
145
146 /**
147  * Given a disk, this function detects, creates and mounts partitions on that disk.
148  * For each detected partition, the filesystem is also detected, and that is mounted if it exists.
149  * If the disk has no partitions, a 'whole disk partition' will be created, the filesystem will
150  * still be detected. This should only be called once per disk, after it has been initialised.
151  *
152  * @param disk The disk to scan for partitions
153  *
154  * @maxirql IRQL_STANDARD
155  */
156
157

```

```

158 void CreateDiskPartitions(struct open_file* disk) {
159     EXACT_IQRL(IQRL_STANDARD);
160
161     struct open_file** partitions = GetPartitionsForDisk(disk);
162
163     if (partitions == NULL || partitions[0] == NULL) {
164         LogWriteSerial("no partition table...\n");
165         struct stat st;
166         int res = VnodeOpStat(disk->node, &st);
167         if (res != 0) {
168             return;
169         }
170
171         struct vnode* whole_disk_partition = CreatePartition(disk, 0, st.st_size, 0, st.st_blksize, 0, false)->node;
172         VnodeOpCreate(disk->node, &whole_disk_partition, GetPartitionNameString(0), 0, 0);
173         return;
174     }
175
176     for (int i = 0; partitions[i]; ++i) {
177         struct vnode* partition = partitions[i]->node;
178         char* str = GetPartitionNameString(i);
179         VnodeOpCreate(disk->node, &partition, str, 0, 0);
180     }
181 }
182
183 void InitDiskPartitionHelper(struct disk_partition_helper* helper) {
184     helper->num_partitions = 0;
185 }
186
187 int DiskFollowHelper(struct disk_partition_helper* helper, struct vnode** out, const char* name) {
188     for (int i = 0; i < helper->num_partitions; ++i) {
189         if (!strcmp(helper->partition_names[i], name)) {
190             *out = helper->partitions[i];
191             return 0;
192         }
193     }
194
195     *out = NULL;
196     return EINVAL;
197 }
198
199 int DiskCreateHelper(struct disk_partition_helper* helper, struct vnode** in, const char* name) {
200     if (helper->num_partitions == MAX_PARTITIONS_PER_DISK) {
201         return EINVAL;
202     }
203
204     helper->partitions[helper->num_partitions] = *in;
205     helper->partition_names[helper->num_partitions] = (char*) name;
206     helper->num_partitions++;
207     return 0;
208 }

```

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
79     inline_memcpy(kernel_addr, user_addr, size);
80     return 0;
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
89     inline_memcpy(user_addr, kernel_addr, size);
90     return 0;
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     }
110
111     } else {
112         int result;
113
114         /*
115          * This is from the kernel's perspective of the operations.
116          */
117         if (untrusted_buffer->direction == TRANSFER_READ) {
118             result = CopyOutOfKernel((const void*) trusted_buffer, untrusted_buffer->address, amount_to_copy);
119         } else {
120             result = CopyIntoKernel(trusted_buffer, (const void*) untrusted_buffer->address, amount_to_copy);
121         }
122     }
123
124     if (result != 0) {
125         return result;
126     }
127 }
128
129 untrusted_buffer->length_remaining -= amount_to_copy;
130 untrusted_buffer->offset += amount_to_copy;
131 untrusted_buffer->address = ((uint8_t*) untrusted_buffer->address) + amount_to_copy;
132
133 return 0;
134 }
135
136 int WriteStringToUsermode(const char* trusted_string, char* untrusted_buffer, uint64_t max_length) {
137     struct transfer tr = CreateTransferWritingToUser(untrusted_buffer, max_length, 0);
138     int result;
139
140     /*
141     * Limit the size of the string by the maximum length. We use <, and a -1 in the other case,
142     * as we need to ensure the null terminator fits.
143     */
144     uint64_t size = strlen(trusted_string) < max_length ? strlen(trusted_string) : max_length - 1;
145     result = PerformTransfer((void*) trusted_string, &tr, size);
146
147     if (result != 0) {
148         return result;
149     }
150
151     uint8_t zero = 0;
152     return PerformTransfer(&zero, &tr, 1);
153 }
154
155 int ReadStringFromUsermode(char* trusted_buffer, const char* untrusted_string, uint64_t max_length) {
156     struct transfer tr = CreateTransferReadingFromUser(untrusted_string, max_length, 0);
157     size_t i = 0;
158
159     while (max_length-- > 1) {
160         char c;
161         int result = PerformTransfer(&c, &tr, 1);
162         if (result != 0) {
163             return result;
164         }
165         trusted_buffer[i++] = c;
166         if (c == 0) {
167             break;
168         }
169     }
170
171     trusted_buffer[i] = 0;
172     return 0;
173 }
174
175 int WriteWordToUsermode(size_t* location, size_t value) {
176     struct transfer io = CreateTransferWritingToUser(location, sizeof(size_t), 0);
177     int res = PerformTransfer(&value, &io, sizeof(size_t));
178     if (io.length_remaining != 0) {
179         return EINVAL;
180     }
181     return res;
182 }
183
184 int ReadWordFromUsermode(size_t* location, size_t* output) {
185     struct transfer io = CreateTransferReadingFromUser(location, sizeof(size_t), 0);
186     int res = PerformTransfer(output, &io, sizeof(size_t));
187     if (io.length_remaining != 0) {
188         return EINVAL;
189     }
190     return res;
191 }
192
193 static struct transfer CreateTransfer(void* addr, uint64_t length, uint64_t offset, int direction, int type) {
194     struct transfer trans;
195     trans.address = addr;
196     trans.direction = direction;
197     trans.length_remaining = length;
198     trans.offset = offset;
199     trans.type = type;
200     return trans;
201 }
202
203

```

```

204 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 int CreateFileDescriptor(struct filedes_table* table, struct open_file* file, int* fd_out, int flags) {
60     if ((flags & O_CLOEXEC) != 0) {
61         return EINVAL;
62     }
63
64     AcquireMutex(table->lock, -1);
65
66     for (int i = 0; i < MAX_FD_PER_PROCESS; ++i) {
67         if (table->entries[i].file == NULL) {
68             table->entries[i].file = file;
69             table->entries[i].flags = flags;
70             ReleaseMutex(table->lock);
71             *fd_out = i;
72             return 0;
73         }
74     }
75
76     ReleaseMutex(table->lock);
77     return EMFILE;
78 }
79
80 int RemoveFileDescriptor(struct filedes_table* table, struct open_file* file) {
81     AcquireMutex(table->lock, -1);
82
83     for (int i = 0; i < MAX_FD_PER_PROCESS; ++i) {
84         if (table->entries[i].file == file) {
85             table->entries[i].file = NULL;
86             ReleaseMutex(table->lock);
87             return 0;
88         }
89     }
90
91     ReleaseMutex(table->lock);
92     return EINVAL;
93 }
94
95 int GetFileFromDescriptor(struct filedes_table* table, int fd, struct open_file** out) {
96     if (out == NULL || fd < 0 || fd >= MAX_FD_PER_PROCESS) {
97         *out = NULL;
98         return out == NULL ? EINVAL : EBADF;
99     }
100
101     AcquireMutex(table->lock, -1);
102     struct open_file* result = table->entries[fd].file;
103     ReleaseMutex(table->lock);
104
105     *out = result;

```

```

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

```

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 <irq.h>
10
11 /*
12  * vfs/vnode.c - Virtual Filesystem Nodes
13  */
14 /* Each vnode represents an abstract file, such as a file, directory or device.
15 */
16
17 /*
18  * Allocate and initialise a vnode. The reference count is initialised to 1.
19 */
20 struct vnode* CreateVnode(struct vnode_operations ops) {
21     struct vnode* node = AllocHeap(sizeof(struct vnode));
22     node->ops = ops;
23     node->reference_count = 1;
24     node->data = NULL;
25     InitSpinlock(&node->reference_count_lock, "vnode refcnt", IRQL_SCHEDULER);
26     return node;
27 }
28
29 /*
30  * Cleanup and free an abstract file node.

```

```

31 */
32 static void DestroyVnode(struct vnode* node) {
33     /*
34      * The lock can't be held during this process, otherwise the lock will
35      * get freed before it is released (which is bad, as we must release it
36      * to get interrupts back on).
37      */
38
39     assert(node != NULL);
40     assert(node->reference_count == 0);
41
42     FreeHeap(node);
43 }
44
45 /*
46  * Ensures that a vnode is valid.
47  */
48 static void CheckVnode(struct vnode* node) {
49     assert(node != NULL);
50
51     if (IsSpinlockHeld(&node->reference_count_lock)) {
52         assert(node->reference_count > 0);
53     } else {
54         AcquireSpinlockIrql(&node->reference_count_lock);
55         assert(node->reference_count > 0);
56         ReleaseSpinlockIrql(&node->reference_count_lock);
57     }
58 }
59
60
61 /*
62  * Increments a vnode's reference counter. Used whenever a vnode is 'given' to someone.
63  */
64 void ReferenceVnode(struct vnode* node) {
65     assert(node != NULL);
66
67     AcquireSpinlockIrql(&node->reference_count_lock);
68     node->reference_count++;
69     ReleaseSpinlockIrql(&node->reference_count_lock);
70 }
71
72 /*
73  * Decrements a vnode's reference counter, destorying it if it reaches zero.
74  * It should be called to free a vnode 'given' to use when it is no longer needed.
75  */
76 void DereferenceVnode(struct vnode* node) {
77     CheckVnode(node);
78     AcquireSpinlockIrql(&node->reference_count_lock);
79
80     assert(node->reference_count > 0);
81     node->reference_count--;
82
83     if (node->reference_count == 0) {
84         VnodeOpClose(node);
85
86         /*
87          * Must release the lock before we delete it so we can put interrupts back on
88          */
89         ReleaseSpinlockIrql(&node->reference_count_lock);
90
91         DestroyVnode(node);
92         return;
93     }
94     ReleaseSpinlockIrql(&node->reference_count_lock);
95 }
96
97
98
99 /*
100  * Wrapper functions for performing operations on a vnode. Also
101  * performs validation on the vnode.
102  */
103 int VnodeOpCheckOpen(struct vnode* node, const char* name, int flags) {
104     if (node == NULL) {
105         return EINVAL;
106     }
107     if (node->ops.check_open == NULL) {
108         return 0;
109     }
110     return node->ops.check_open(node, name, flags);
111 }
112
113 int VnodeOpRead(struct vnode* node, struct transfer* io) {
114     if (node == NULL || node->ops.read == NULL || io->direction != TRANSFER_READ) {
115         return EINVAL;
116     }
117     return node->ops.read(node, io);
118 }
119
120 int VnodeOpWrite(struct vnode* node, struct transfer* io) {
121     if (node == NULL || node->ops.write == NULL || io->direction != TRANSFER_WRITE) {
122         return EINVAL;
123     }
124     return node->ops.write(node, io);
125 }
126
127 int VnodeOpIoctl(struct vnode* node, int command, void* buffer) {
128     if (node == NULL || node->ops.ioctl == NULL) {
129         return EINVAL;
130     }
131     return node->ops.ioctl(node, command, buffer);
132 }
133
134 bool VnodeOpIsSeekable(struct vnode* node) {
135     if (node == NULL || node->ops.is_seekable == NULL) {
136         return false;
137     }
138     return node->ops.is_seekable(node);
139 }
140
141 int VnodeOpCheckTty(struct vnode* node) {
142     if (node == NULL || node->ops.check_tty == NULL) {
143         return ENOTTY;
144     }
145     return node->ops.check_tty(node);
146 }
147
148 int VnodeOpClose(struct vnode* node) {
149     if (node == NULL || node->reference_count != 0) {
150         return EINVAL;
151     }
152     if (node->ops.close == NULL) {
153         return 0;
154     }
155     return node->ops.close(node);
156 }
157
158 int VnodeOpCreate(struct vnode* node, struct vnode** out, const char* name, int flags, mode_t mode) {
159     if (node == NULL || node->ops.create == NULL) {
160         return EINVAL;
161     }

```



```

161     }
162     return node->ops.create(node, out, name, flags, mode);
163 }
164
165 int VnodeOpTruncate(struct vnode* node, off_t offset) {
166     if (node == NULL || node->ops.truncate == NULL) {
167         return EINVAL;
168     }
169     return node->ops.truncate(node, offset);
170 }
171
172 int VnodeOpFollow(struct vnode* node, struct vnode** new_node, const char* name) {
173     if (node == NULL || node->ops.follow == NULL) {
174         return ENOTDIR;
175     }
176     return node->ops.follow(node, new_node, name);
177 }
178
179 uint8_t VnodeOpDirentType(struct vnode* node) {
180     if (node == NULL || node->ops.dirent_type == NULL) {
181         return DT_UNKNOWN;
182     }
183     return node->ops.dirent_type(node);
184 }
185
186 int VnodeOpStat(struct vnode* node, struct stat* stat) {
187     if (node == NULL || node->ops.stat == NULL) {
188         return EINVAL;
189     }
190     return node->ops.stat(node, stat);
191 }
192
193 int VnodeOpWait(struct vnode* node, int flags, uint64_t timeout_ms) {
194     if (node == NULL) {
195         return EINVAL;
196     }
197     if (node->ops.wait == NULL) {
198         return 0;
199     }
200     return node->ops.wait(node, flags, timeout_ms);
201 }

```

File: ./vfs/vfs.c

```

1  #include <vfs.h>
2  #include <spinlock.h>
3  #include <irql.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 /*
22  * Maximum length of a component of a filepath (e.g. an file/directory's individual name).
23  */
24 #define MAX_COMPONENT_LENGTH 128
25
26 /*
27  * Maximum total length of a path.
28  */
29 #define MAX_PATH_LENGTH 2000
30
31 /*
32  * Maximum number of symbolic links to dereference in a path before returning ELOOP.
33  */
34 #define MAX_LOOP 5
35
36
37 /*
38  * A structure for mounted devices and filesystems.
39  */
40 struct mounted_file {
41     /* The vnode representing the device / root directory of a filesystem */
42     struct open_file* node;
43
44     /* What the device / filesystem mount is called */
45     char* name;
46 };
47
48 static struct spinlock vfs_lock;
49 static struct avl_tree* mount_points = NULL;
50
51 int MountedDeviceComparator(void* a_, void* b_) {
52     struct mounted_file* a = a_;
53     struct mounted_file* b = b_;
54     return strcmp(a->name, b->name);
55 }
56
57 void InitVfs(void) {
58     InitSpinlock(&vfs_lock, "vfs", IRQL_SCHEDULER);
59     mount_points = AvlTreeCreate();
60     AvlTreeSetComparator(mount_points, MountedDeviceComparator);
61 }
62
63 static int CheckValidComponentName(const char* name)
64 {
65     assert(name != NULL);
66
67     if (name[0] == 0) {
68         return EINVAL;
69     }
70
71     for (int i = 0; name[i]; ++i) {
72         char c = name[i];
73
74         if (c == '/' || c == '\\\ ' || c == ':') {
75             return EINVAL;
76         }
77     }
78
79     return 0;
80 }
81
82 static int DoesMountPointExist(const char* name) {
83     assert(name != NULL);
84     assert(IsSpinlockHeld(&vfs_lock));

```

```

85
86     struct mounted_file target;
87     target.name = (char*) name;
88     if (AvlTreeContains(mount_points, &target)) {
89         return EEXIST;
90     }
91
92     return 0;
93 }
94
95 /*
96  * Given a filepath, and a pointer to an index within that filepath (representing where
97  * start searching), copies the next component into an output buffer of a given length.
98  * The index is updated to point to the start of the next component, ready for the next call.
99  *
100  * This also handles duplicated and trailing forward slashes.
101  */
102 static int GetPathComponent(const char* path, int* ptr, char* output_buffer, int max_output_length, char delimiter) {
103     int i = 0;
104
105     assert(path != NULL);
106     assert(ptr != NULL);
107     assert(max_output_length >= 1);
108     assert(delimiter != 0);
109     assert(output_buffer != NULL);
110     assert(strlen(path) >= 1);
111     assert(*ptr >= 0 && *ptr < (int) strlen(path));
112
113     /*
114      * These were meant to be caught at a higher level, so we can apply the current
115      * working directory or the current drive.
116      */
117     assert(path[0] != '/');
118     assert(path[0] != ':');
119
120     output_buffer[0] = 0;
121
122     while (path[*ptr] && path[*ptr] != delimiter) {
123         if (i >= max_output_length - 1) {
124             return ENAMETOOLONG;
125         }
126
127         /*
128          * Ensure we always have a null terminated string.
129          */
130         output_buffer[i++] = path[*ptr];
131         output_buffer[i] = 0;
132         (*ptr)++;
133     }
134
135     /*
136      * Skip past the delimiter (unless we are at the end of the string),
137      * as well as any trailing slashes (which could be after a slash delimiter, or
138      * after a colon).
139      */
140     if (path[*ptr]) {
141         do {
142             (*ptr)++;
143         } while (path[*ptr] == '/');
144     }
145
146     /*
147      * Ensure that there are no colons or backslashes in the filename itself.
148      */
149     return CheckValidComponentName(output_buffer);
150 }
151
152 static int GetFinalPathComponent(const char* path, char* output_buffer, int max_output_length) {
153     int path_ptr = 0;
154
155     int status = GetPathComponent(path, &path_ptr, output_buffer, max_output_length, ':');
156     if (status) {
157         return status;
158     }
159
160     while (path_ptr < (int) strlen(path)) {
161         status = GetPathComponent(path, &path_ptr, output_buffer, max_output_length, '/');
162         if (status) {
163             return status;
164         }
165     }
166
167     return 0;
168 }
169
170 /*
171  * Takes in a device name, without the colon, and returns its vnode.
172  * If no such device is mounted, it returns NULL.
173  */
174 static struct open_file* GetMountFromName(const char* name) {
175     assert(name != NULL);
176
177     if (mount_points == NULL) {
178         return NULL;
179     }
180
181     struct mounted_file target;
182     target.name = (char*) name;
183     struct mounted_file* mount = AvlTreeGet(mount_points, (void*) &target);
184     return mount->node;
185 }
186
187 int PAGEABLE_CODE_SECTION AddVfsMount(struct vnode* node, const char* name) {
188     MAX_IrqL[IRQL_PAGE_FAULT];
189
190     if (name == NULL || node == NULL) {
191         return EINVAL;
192     }
193
194     if (strlen(name) >= MAX_COMPONENT_LENGTH) {
195         return ENAMETOOLONG;
196     }
197
198     int status = CheckValidComponentName(name);
199     if (status != 0) {
200         return status;
201     }
202
203     AcquireSpinlockIrql(&vfs_lock);
204
205     if (DoesMountPointExist(name) == EEXIST) {
206         ReleaseSpinlockIrql(&vfs_lock);
207         return EEXIST;
208     }
209
210     struct mounted_file* mount = AllocHeap(sizeof(struct mounted_file));
211     mount->name = strdup(name);
212     mount->node = CreateOpenFile(node, 0, 0, true, true);
213
214     AvlTreeInsert(mount_points, (void*) mount);

```

```

215
216 LogWriteSerial("MOUNTED TO THE VFS: %s\n", name);
217
218     ReleaseSpinlockIrql(&vfs_lock);
219     return 0;
220
221
222 int PAGEABLE_CODE_SECTION RemoveVfsMount(const char* name) {
223     MAX_IRQL IRQL_PAGE_FAULT;
224
225     if (name == NULL) {
226         return EINVAL;
227     }
228
229     if (CheckValidComponentName(name) != 0) {
230         return EINVAL;
231     }
232
233     AcquireSpinlockIrql(&vfs_lock);
234
235     /*
236      * Scan through the mount table for the device
237      */
238     struct mounted_file target;
239     target.name = (char*) name;
240
241     struct mounted_file* actual = AvlTreeGet(mount_points, &target);
242     if (actual == NULL) {
243         ReleaseSpinlockIrql(&vfs_lock);
244         return ENODEV;
245     }
246
247     assert(!strcmp(actual->name, name));
248
249     /*
250      * Decrement the reference that was initially created way back in
251      * vfs_add_device in the call to dev_create_vnode (the vnode dereference),
252      * and then the open file that was created alongside it.
253      */
254     DereferenceVnode(actual->node->node);
255     DereferenceOpenFile(actual->node);
256
257     AvlTreeDelete(mount_points, actual);
258     FreeHeap(actual->name);
259
260     ReleaseSpinlockIrql(&vfs_lock);
261     return 0;
262 }
263
264 static void CleanupVnodeStack(struct stack_adt* stack) {
265     /*
266      * We need to call dereference on each vnode in the stack before we
267      * can call StackAdtDestroy.
268      */
269     while (StackAdtSize(stack) > 0) {
270         struct vnode* node = StackAdtPop(stack);
271         DereferenceVnode(node);
272     }
273
274     StackAdtDestroy(stack);
275 }
276
277 /*
278  * Given an absolute filepath, returns the vnode representing
279  * the file, directory or device.
280  *
281  * Should only be called by vfs_open as the reference count will be incremented.
282  */
283 static int GetVnodeFromPath(const char* path, struct vnode** out, bool want_parent) {
284     assert(path != NULL);
285     assert(out != NULL);
286
287     LogWriteSerial("GetVnodeFromPath: path %s\n", path);
288
289     if (strlen(path) == 0) {
290         return EINVAL;
291     }
292     if (strlen(path) > MAX_PATH_LENGTH) {
293         return ENAMETOOLONG;
294     }
295
296     int path_ptr = 0;
297     char component_buffer[MAX_COMPONENT_LENGTH];
298
299     int err = GetPathComponent(path, &path_ptr, component_buffer, MAX_COMPONENT_LENGTH, '/');
300     if (err != 0) {
301         return err;
302     }
303
304     struct open_file* current_file = GetMountFromName(component_buffer);
305
306     /*
307      * No root device found, so we can't continue.
308      */
309     if (current_file == NULL || current_file->node == NULL) {
310         return ENODEV;
311     }
312
313     struct vnode* current_vnode = current_file->node;
314
315     /*
316      * This will be dereferenced either as we go through the loop, or
317      * after a call to vfs_close (this function should only be called
318      * by vfs_open).
319      */
320     ReferenceVnode(current_vnode);
321
322     char component[MAX_COMPONENT_LENGTH + 1];
323
324     /*
325      * To go back to a parent directory, we need to keep track of the previous component.
326      * As we can go back through many parents, we must keep track of all of them, hence we
327      * use a stack to store each vnode we encounter. We will not dereference the vnodes
328      * on the stack until the end using cleanup_vnode_stack.
329      */
330     struct stack_adt* previous_components = StackAdtCreate();
331
332     /*
333      * Iterate over the rest of the path.
334      */
335     while (path_ptr < (int) strlen(path)) {
336         int status = GetPathComponent(path, &path_ptr, component, MAX_COMPONENT_LENGTH, '/');
337         if (status != 0) {
338             DereferenceVnode(current_vnode);
339             CleanupVnodeStack(previous_components);
340             return status;
341         }
342
343         if (!strcmp(component, ".") || !strcmp(component, "..")) {
344             /*

```

```

345  * This doesn't change where we point to.
346  */
347  continue;
348
349  } else if (!strcmp(component, "..")) {
350      if (StackAdtSize(previous_components) == 0) {
351          /*
352           * We have reached the root. Going 'further back' than the root
353           * on Linux just keeps us at the root, so don't do anything here.
354           */
355
356      } else {
357          /*
358           * Pop the previous component and use it.
359           */
360          current_vnode = StackAdtPop(previous_components);
361      }
362
363      continue;
364  }
365
366  /*
367   * Use a separate pointer so that both inputs don't point to the same
368   * location. vnode follow either increments the reference count or creates
369   * a new vnode with a count of one.
370   */
371  struct vnode* next_vnode = NULL;
372  status = VnodeOpFollow(current_vnode, &next_vnode, component);
373  if (status != 0) {
374      DereferenceVnode(current_vnode);
375      CleanupVnodeStack(previous_components);
376      return status;
377  }
378
379  /*
380   * We have a component that can be backtracked to, so add it to the stack.
381   *
382   * Also note that vnode follow adds a reference count, so current_vnode
383   * needs to be dereferenced. Conveniently, all components that need to be
384   * put on the stack also need dereferencing, and vice versa.
385   *
386   * The final vnode we find will not be added to the stack and dereferenced
387   * as we won't get here.
388   */
389  StackAdtPush(previous_components, current_vnode);
390  current_vnode = next_vnode;
391  }
392
393  int status = 0;
394
395  if (want_parent) {
396      /*
397       * Operations that require us to get the parent don't work if we are already
398       * at the root.
399       */
400      if (StackAdtSize(previous_components) == 0) {
401          status = EINVAL;
402      } else {
403          *out = StackAdtPop(previous_components);
404      }
405  } else {
406      *out = current_vnode;
407  }
408
409  CleanupVnodeStack(previous_components);
410  return status;
411
412
413
414 int OpenFile(const char* path, int flags, mode_t mode, struct open_file** out) {
415     EXACT_IOCTL(IRQL_STANDARD);
416
417     if (path == NULL || out == NULL || strlen(path) <= 0) {
418         return EINVAL;
419     }
420
421     int status;
422     char name[MAX_COMPONENT_LENGTH + 1];
423     status = GetFinalPathComponent(path, name, MAX_COMPONENT_LENGTH);
424     if (status) {
425         return status;
426     }
427
428     /*
429      * Grab the vnode from the path.
430      */
431     struct vnode* node;
432
433     /*
434      * Lookup a (hopefully) existing file.
435      */
436     status = GetVnodeFromPath(path, &node, false);
437
438     if (flags & O_CREAT) {
439         if (status == ENOENT) {
440             /*
441              * Get the parent folder.
442              */
443             status = GetVnodeFromPath(path, &node, true);
444             if (status)
445                 return status;
446
447             struct vnode* child;
448             status = VnodeOpCreate(node, &child, name, flags, mode);
449             DereferenceVnode(node);
450
451             if (status) {
452                 return status;
453             }
454
455             node = child;
456
457         } else if (flags & O_EXCL) {
458             /*
459              * The file already exists (as we didn't get ENOENT), but we were passed O_EXCL so we
460              * must give an error. If O_EXCL isn't passed, then O_CREAT will just open the existing file.
461              */
462             return EEXIST;
463         }
464     }
465
466     if (status != 0) {
467         return status;
468     }
469
470     status = VnodeOpCheckOpen(node, name, flags & (O_ACCMODE | O_NONBLOCK));
471     if (status) {
472         DereferenceVnode(node);
473         return status;
474     }

```

```

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

```

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         DbgScreenPutchchar(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
38 codepoints[out++] = codepoint;
39 }
40
41 *out_length = out;
42 return 0;
43 }
44
45 int Utf8ToCodepoints(uint8_t* utf8, uint32_t* codepoints, int in_length, int* out_length) {
46     int in = 0;
47     int out = 0;
48
49     while (in < in_length) {
50         if (out == out_length) {
51             return ENAMETOOLONG;
52         }
53
54         uint32_t codepoint = utf8[in++];
55         if ((codepoint >> 3) == 0xE) {
56             if (in + 2 >= in_length) {
57                 return EINVAL;
58             }
59             codepoint &= 0x7;
60             codepoint <<= 6;
61             uint8_t next = utf8[in++];
62             if ((next >> 6) != 0x2) {
63                 return EINVAL;
64             }
65             codepoint |= next & 0x3F;
66             codepoint <<= 6;
67             next = utf8[in++];
68             if ((next >> 6) != 0x2) {
69                 return EINVAL;
70             }
71             codepoint |= next & 0x3F;
72             next = utf8[in++];
73             if ((next >> 6) != 0x2) {
74                 return EINVAL;
75             }
76             codepoint |= next & 0x3F;
77         } else if ((codepoint >> 4) == 0xE) {
78             if (in + 1 >= in_length) {
79                 return EINVAL;
80             }
81             codepoint &= 0xF;
82             codepoint <<= 6;
83             uint8_t next = utf8[in++];
84             if ((next >> 6) != 0x2) {
85                 return EINVAL;
86             }
87             codepoint |= next & 0x3F;
88             codepoint <<= 6;
89             next = utf8[in++];
90             if ((next >> 6) != 0x2) {
91                 return EINVAL;
92             }
93             codepoint |= next & 0x3F;
94         } else if ((codepoint >> 5) == 0x6) {
95             if (in >= in_length) {
96                 return EINVAL;
97             }
98             codepoint &= 0x1F;
99             codepoint <<= 6;
100             uint8_t next = utf8[in++];
101             if ((next >> 6) != 0x2) {
102                 return EINVAL;
103             }
104             codepoint |= next & 0x3F;
105         }
106
107         if (codepoint >= 0x80) {
108             return EINVAL;
109         }
110
111         if (codepoint >= 0xD800 && codepoint <= 0xDFFF) {
112             return EINVAL;
113         }
114
115         codepoints[out++] = codepoint;
116     }
117
118     *out_length = out;
119     return 0;
120 }
121
122 int CodepointsToUtf16(uint32_t* codepoints, uint16_t* utf16, int in_length, int* out_length) {
123     int in = 0;
124     int out = 0;
125     while (in < in_length) {
126         if (out == out_length) {
127             return ENAMETOOLONG;
128         }
129
130         uint32_t codepoint = codepoints[in++];
131
132         if (codepoint >= 0xD800 && codepoint <= 0xDFFF) {
133             return EINVAL;
134         }
135
136         if (codepoint > 0x10FFFF) {
137             return EINVAL;
138         }
139
140         if (codepoint >= 0x10000) {
141             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     }
150 } else {
151     utf16[out++] = codepoint;
152 }
153 }
154
155 *out_length = out;
156 return 0;
157 }
158
159 int CodepointsToUtf8(uint32_t* codepoints, uint8_t* utf8, int in_length, int* out_length) {
160     int in = 0;
161     int out = 0;
162     while (in < in_length) {
163         if (out == *out_length) {
164             return ENAMETOOLONG;
165         }
166
167         uint32_t codepoint = codepoints[in++];
168
169         if (codepoint >= 0xD800 && codepoint <= 0xDFFF) {
170             return EINVAL;
171         }
172
173         if (codepoint <= 0x7F) {
174             utf8[out++] = codepoint;
175         } else if (codepoint <= 0x7FF) {
176             if (out + 2 > *out_length) {
177                 return ENAMETOOLONG;
178             }
179             utf8[out++] = 0xC0 | (codepoint >> 6);
180             utf8[out++] = 0x80 | (codepoint & 0x3F);
181         } else if (codepoint <= 0xFFFF) {
182             if (out + 3 > *out_length) {
183                 return ENAMETOOLONG;
184             }
185             utf8[out++] = 0xE0 | (codepoint >> 12);
186             utf8[out++] = 0x80 | ((codepoint >> 6) & 0x3F);
187             utf8[out++] = 0x80 | (codepoint & 0x3F);
188         } else if (codepoint <= 0x10FFF) {
189             if (out + 4 > *out_length) {
190                 return ENAMETOOLONG;
191             }
192             utf8[out++] = 0xF0 | (codepoint >> 18);
193             utf8[out++] = 0x80 | ((codepoint >> 12) & 0x3F);
194             utf8[out++] = 0x80 | ((codepoint >> 6) & 0x3F);
195             utf8[out++] = 0x80 | (codepoint & 0x3F);
196         } else {
197             return EINVAL;
198         }
199     }
200 }
201
202 *out_length = out;
203 return 0;
204 }
205
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 DeferPutchar(void* v) {
20     video_driver.putchar((char) (size_t) v);
21 }
22
23 void DbgScreenPutchar(char c) {
24     if (video_driver.putchar != NULL) {
25         DeferUntilIrql(IRQL_STANDARD, DeferPutchar, (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_MALLOC_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;
60     dummy.filename = (char*) name;
61
62     AvlTreeSetComparator(loaded_drivers, DriverTableComparatorByName);
63     struct loaded_driver* res = AvlTreeGet(loaded_drivers, &dummy);
64     if (res == NULL) {
65         return 0;
66     }
67     return res->relocation_point;
68 }
69
70 static struct loaded_driver* GetDriverFromAddress(size_t relocation_point) {
71     AcquireMutex(driver_table_lock, -1);
72
73     struct loaded_driver dummy;
74     dummy.relocation_point = relocation_point;
75
76     AvlTreeSetComparator(loaded_drivers, DriverTableComparatorByRelocationPoint);
77     struct loaded_driver* res = AvlTreeGet(loaded_drivers, &dummy);
78     ReleaseMutex(driver_table_lock);
79     return res;
80 }
81
82 size_t GetDriverAddress(const char* name) {
83     EXACT_IRQL(IRQL_STANDARD);
84
85     AcquireMutex(driver_table_lock, -1);
86     size_t res = GetDriverAddressWithLockHeld(name);
87     ReleaseMutex(driver_table_lock);
88     return res;
89 }
90
91 void InitSymbolTable(void) {
92     driver_table_lock = CreateMutex("drv table");
93     symbol_table_lock = CreateMutex("sym table");
94
95     loaded_drivers = AvlTreeCreate();
96     symbol_table = AvlTreeCreate();
97     AvlTreeSetComparator(symbol_table, SymbolComparator);
98
99     struct open_file* kernel_file;
100     if (OpenFile("sys:/kernel.exe", O_RDONLY, 0, &kernel_file)) {
101         Panic(PANIC_NO_FILESYSTEM);
102     }
103     ArchLoadSymbols(kernel_file, 0);
104     CloseFile(kernel_file);
105 }
106
107 static bool DoesSymbolContainIllegalCharacters(const char* symbol) {
108     for (int i = 0; symbol[i]; ++i) {
109         if (!isalnum(symbol[i]) && symbol[i] != '_' ) {
110             return true;
111         }
112     }
113     return strlen(symbol) == 0;
114 }
115
116 void AddSymbol(const char* symbol, size_t address) {
117     EXACT_IRQL(IRQL_STANDARD);
118
119     if (DoesSymbolContainIllegalCharacters(symbol)) {
120         return;
121     }
122
123     struct symbol* entry = AllocHeap(sizeof(struct symbol));
124     entry->name = strdup(symbol);
125     entry->addr = address;
126
127     AcquireMutex(symbol_table_lock, -1);
128     if (AvlTreeContains(symbol_table, entry)) {
129         /*
130          * The kernel has some symbols declared 'static' to file scope, with
131          * duplicate names (e.g. in /dev each file has its own 'Stat'). These
132          * get exported for some reason so we end up with duplicate names. We
133          * must ignore these to avoid AVL issues. They are safe to ignore, as
134          * they were meant to be 'static' anyway.
135          */
136         /*
137          * TODO: there may be issues if device drivers try to create their own
138          * methods with those names (?) e.g. they use the standard template
139          * and have their own 'Stat'.
140          */
141         FreeHeap(entry);
142     } else {
143         AvlTreeInsert(symbol_table, entry);
144     }
145     ReleaseMutex(symbol_table_lock);
146 }
147
148 /*
149 * Do not name a (global) function pointer you receive from this the same as the
150 * actual function - this may cause issues with duplicate symbols.
151 * e.g. don't do this at global scope:
152 *
153 * void (*MyFunc)(void) = GetSymbolAddress("MyFunc");
154 *
155 * Instead, do void (*_MyFunc)(void) or something.
156 */
157 size_t GetSymbolAddress(const char* symbol) {
158     EXACT_IRQL(IRQL_STANDARD);
159
160     struct symbol dummy;
161     dummy.name = symbol;
162
163     AcquireMutex(symbol_table_lock, -1);
164     struct symbol* result = AvlTreeGet(symbol_table, &dummy);
165     ReleaseMutex(symbol_table_lock);
166
167     if (result == NULL) {
168         return 0;
169     } else {
170         assert(strcmp(result->name, symbol));
171         return result->addr;
172     }
173 }
174
175
176 static int LoadDriver(const char* name) {
177     struct open_file* file;
178     int res;
179     if ((res = OpenFile(name, O_RDONLY, 0, &file))) {
180         return res;
181     }
182
183     struct loaded_driver* drv = AllocHeap(sizeof(struct loaded_driver));
184     drv->filename = strdup_pageable(name);
185     drv->quick_relocation_table = NULL;
186     if ((res = ArchLoadDriver(&drv->relocation_point, file, &drv->quick_relocation_table))) {

```

```

187     return res;
188 }
189
190 assert(drv->quick_relocation_table != NULL);
191
192 AvlTreeInsert(loaded_drivers, drv);
193 ArchLoadSymbols(file, drv->relocation_point - 0xD0000000); // TODO: ??? GET RID OF ARCH SPECIFIC DETAILS (0xD0000000)
194 return 0;
195
196
197 int RequireDriver(const char* name) {
198     EXACT_IROL(IROL_STANDARD);
199
200     LogWriteSerial("Requiring driver: %s\n", name);
201
202     AcquireMutex(driver_table_lock, -1);
203
204     if (GetDriverAddressWithLockHeld(name) != 0) {
205         ReleaseMutex(driver_table_lock);
206
207         /*
208          * Not an error that it's already loaded - ideally no one should care if it has already been loaded
209          * or not. Hence we give 0 (success), not EALREADY.
210          */
211         return 0;
212     }
213
214     int res = LoadDriver(name);
215     ReleaseMutex(driver_table_lock);
216     return res;
217 }
218
219 void SortQuickRelocationTable(struct quick_relocation_table* table) {
220     struct quick_relocation* entries = table->entries;
221     qsort_pageable((void*) entries, table->used_entries, sizeof(struct quick_relocation), QuickRelocationTableComparator);
222 }
223
224 void AddToQuickRelocationTable(struct quick_relocation_table* table, size_t addr, size_t val) {
225     assert(table->used_entries < table->total_entries);
226     table->entries[table->used_entries].address = addr;
227     table->entries[table->used_entries].value = val;
228     table->used_entries++;
229 }
230
231 struct quick_relocation_table* CreateQuickRelocationTable(int count) {
232     struct quick_relocation_table* table = AllocHeap(sizeof(struct quick_relocation_table));
233     struct quick_relocation* entries = (struct quick_relocation*) MapVirt(0, 0, count * sizeof(struct quick_relocation), VM_READ | VM_WRITE, NULL, 0);
234     table->entries = entries;
235     table->used_entries = 0;
236     table->total_entries = count;
237     return table;
238 }
239
240 static int BinarySearchComparator(const void* a, const void* b) {
241     struct quick_relocation a = *((struct quick_relocation*) a);
242     struct quick_relocation b = *((struct quick_relocation*) b);
243
244     size_t page_a = a.address / ARCH_PAGE_SIZE;
245     size_t page_b = b.address / ARCH_PAGE_SIZE;
246
247     if (page_a > page_b) {
248         return 1;
249     } else if (page_a < page_b) {
250         return -1;
251     } else {
252         return 0;
253     }
254 }
255
256
257
258 static void ApplyRelocationsToPage(struct quick_relocation_table* table, size_t virtual) {
259     struct quick_relocation target;
260     target.address = virtual;
261     LogWriteSerial("ApplyRelocationsToPage: looking for 0x%X\n", virtual);
262
263     struct quick_relocation* entry = bsearch(&target, table->entries, table->used_entries, sizeof(struct quick_relocation), BinarySearchComparator);
264     if (entry == NULL) {
265         PanicEx(PANIC_ASSERTION_FAILURE, "quick relocation table doesn't contain lookup - bsearch or qsort is probably bugged");
266     }
267
268     /*
269      * As we only did a binary search to look for anything on that page, we might be halfway through the page!
270      * Move back through the entries until we find the start of page (or the first entry in the table).
271      */
272     while ((entry->address / ARCH_PAGE_SIZE == virtual / ARCH_PAGE_SIZE
273            || (entry->address - sizeof(size_t) + 1) == virtual / ARCH_PAGE_SIZE) && entry != table->entries) {
274         entry -= 1;
275     }
276
277     /*
278      * We went past the last one, so need to forward onto it again - unless we hit the start of the table.
279      */
280     if (entry != table->entries) {
281         entry += 1;
282     }
283
284     /*
285      * We also need to lock these, as (and YES this has actually happened before):
286      * - if the relocation is on the boundary of the next one, and the next one is not present
287      * - and the next one is relocatable
288      * Then:
289      * - we make the next page writable
290      * - we try to do the straddle relocation
291      * - that causes a fault on the next page (not present)
292      * - that one is also relocatable, so it enables writing on that page
293      * - it finishes, and unmarks it as writable
294      * - we fail to do the relocation, as it is no longer writable
295      *
296      * By locking it first, we force the relocations on the second page to happen first, and then we can
297      * mark it as writable.
298      */
299
300     bool needs_write_low = (GetVirtPermissions(virtual) & VM_WRITE) == 0;
301     bool needs_write_high = false;
302     bool need_unlock_high = false;
303
304     if (needs_write_low) {
305         SetVirtPermissions(virtual, VM_WRITE, 0);
306     }
307
308     size_t final_address = table->entries[table->used_entries - 1].address;
309
310     while (entry->address / ARCH_PAGE_SIZE == virtual / ARCH_PAGE_SIZE
311            || (entry->address - sizeof(size_t) + 1) == virtual / ARCH_PAGE_SIZE) {
312
313         if ((entry->address + sizeof(size_t) + 1) / ARCH_PAGE_SIZE != virtual / ARCH_PAGE_SIZE) {
314             need_unlock_high = !LockVirt(virtual + ARCH_PAGE_SIZE);
315             needs_write_high = (GetVirtPermissions(virtual + ARCH_PAGE_SIZE) & VM_WRITE) == 0;
316

```

```

317         if (needs_write_high) {
318             SetVirtPermissions(virtual + ARCH_PAGE_SIZE, VM_WRITE, 0);
319         }
320     }
321
322     size_t* ref = (size_t*) entry->address;
323     *ref = entry->value;
324
325     if (entry->address == final_address) {
326         break;
327     }
328     entry += 1;
329 }
330
331 if (needs_write_low) {
332     SetVirtPermissions(virtual, 0, VM_WRITE);
333 }
334 if (needs_write_high) {
335     SetVirtPermissions(virtual + ARCH_PAGE_SIZE, 0, VM_WRITE);
336 }
337 if (need_unlock_high) {
338     UnlockVirt(virtual + ARCH_PAGE_SIZE);
339 }
340 }
341
342 void PerformDriverRelocationOnPage(struct vas*, size_t relocation_base, size_t virt) {
343     LogWriteSerial("PerformDriverRelocationOnPage A\n");
344     struct loaded_driver* drv = GetDriverFromAddress(relocation_base);
345     if (drv == NULL) {
346         PanicEx(PANIC_ASSERTION_FAILURE, "PerformDriverRelocationOnPage");
347     }
348     LogWriteSerial("PerformDriverRelocationOnPage B. driver at 0x%X\n", drv);
349     ApplyRelocationsToPage(drv->quick_relocation_table, virt);
350 }
351 }

```

File: /Makefile

```

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

LINK_FLAGS = -fuse-lld=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 =

OBJECTS = $(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))

oskernel: $(OBJECTS) $(ASMOBJECTS) $(HOBECTS)
$(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);
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, mode) DestroySemaphore(mtx, mode)

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.
*
* is_seekable: default EINVAL
* Returns true if seek can be called on the file.
*
* check_tty: default ENOTTY
* Returns 0 if a terminal, or ENOTTY otherwise.
*
* 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).
*
* dirent_type: default DT_UNKNOWN
* Returns the type of file, either DT_DIR, DT_REG, DT_BLK, DT_CHR, DT_FIFO, DT_LNK, or DT_UNKNOWN.
*
*/

```

```

#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)

```

```

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 (*stat)(struct vnode* node, struct stat* st);
int (*check_tty)(struct vnode* node);
int (*wait)(struct vnode* node, int flags, uint64_t timeout_ms);

```

```

bool (*is_seekable)(struct vnode* node);
uint8_t (*dirent_type)(struct vnode* node);
};

```

```

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

```

```

/*
* Allocates a new vnode for a given set of operations.
*/

```

```

struct vnode* CreateVnode(struct vnode_operations ops);
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);
bool VnodeOpIsSeekable(struct vnode* node);
int VnodeOpCheckTty(struct vnode* node);
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 VnodeOpStat(struct vnode* node, struct stat* st);
int VnodeOpWait(struct vnode* node, int flags, uint64_t timeout_ms);

```

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_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 GetFileSize(struct open_file* file, off_t* size);

```

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))

#ifdef 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_irq, platform_irq_context_t* context);
void UnhandledFault(void);

```

File: `/include/process.h`

```

#pragma once

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

struct filedes_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* prcss);

struct filedes_table* GetFileDescriptorTable(struct process* prcss);

void AddThreadToProcess(struct process* prcss, 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);
```

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

#ifndef 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 IRQ. 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_irql, platform_irq_context_t* context) {
30     int irql = RaiseIrql(required_irql);
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             iter = LinkedListGetNextNode(iter);
47         }
48     }
49
50     LowerIrql(irql);
51 }
52
53 void UnhandledFault(void) {
54     if (GetProcess() != NULL) {
55         LogWriteSerial("unhandled fault...\n");
56         TerminateThread(GetThread());
57     } else {
58         Panic(PANIC_UNHANDLED_KERNEL_EXCEPTION);
59     }
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 = IRQL_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", IRQL_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_IRQL(IRQL_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/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 <filesdes.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 <filesdes.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         struct stat st;
37         if ((res = VnodeOpStat(file->node, &st)) {
38             return res;
39         }
40         offset += st.st_size;
41     }
42     else if (whence != SEEK_SET) {
43         return EINVAL;
44     }
45
46     file->seek_position = offset;
47
48     io = CreateTransferWritingToUser((void*) pos_ptr, sizeof(off_t), 0);
49     return PerformTransfer(&offset, &io, sizeof(off_t));
50 }
```

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 <fcntl.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         struct stat st;
23         if ((res = VnodeOpStat(file->node, &st)) {
24             return res;
25         }
26
27         file->seek_position = st.st_size;
28     }
29
30     struct transfer io = CreateTransferWritingToUser((uint8_t*) buffer, size, file->seek_position);
31     io.direction = write ? TRANSFER_WRITE : TRANSFER_READ;
32
33     if ((res = (write ? WriteFile : ReadFile)(file, &io)) {
34         return res;
35     }
36
37     size_t br = size - io.length_remaining;
38     file->seek_position += br;
39
40     return WriteWordToUsermode((size_t*) br_out, br);
41 }

```

File: ./sys/calls/terminate.c

```

1 #include <syscall.h>
2 #include <errno.h>
3 #include <_syscallnum.h>
4 #include <thread.h>
5
6 int SysTerminate(size_t, size_t, size_t, size_t, size_t) {
7     TerminateThread(GetThread());
8     return EFAULT;
9 }

```

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     }
27     else if (dup_num == 2) {
28         return DuplicateFileDescriptor2(table, old_fd, new_fd, flags);
29     }
30     else {
31         return EINVAL;
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 };
19
20 static const system_call_t system_call_table[_SYSCALL_NUM_ENTRIES] = {
21     [SYSCALL_YIELD] = SysYield,
22     [SYSCALL_TERMINATE] = SysTerminate,
23     [SYSCALL_MAPVIRT] = SysMapVirt,
24     [SYSCALL_UNMAPVIRT] = SysUnmapVirt,
25     [SYSCALL_OPEN] = SysOpen,
26     [SYSCALL_READWRITE] = SysReadWrite,
27     [SYSCALL_CLOSE] = SysClose,
28     [SYSCALL_SEEK] = SysSeek,
29     [SYSCALL_DUP] = SysDup,
30 };
31
32 int HandleSystemCall(int call, size_t a, size_t b, size_t c, size_t d, size_t e) {
33     if (call >= _SYSCALL_NUM_ENTRIES) {
34         return ENOSYS;
35     }
36     return system_call_table[call](a, b, c, d, e);
37 }
38 }

```

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 <irq.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     } else if (num == ISR_PAGE_FAULT) {
33         extern size_t x86GetCr2();
34
35         int type = 0;
36         if (r->err_code & 1) {
37             type |= VM_READ;
38         }
39         if (r->err_code & 2) {
40             type |= VM_WRITE;
41         }
42         if (r->err_code & 4) {
43             type |= VM_USER;
44         }
45         if (r->err_code & 16) {
46             type |= VM_EXEC;
47         }
48
49         LogWriteSerial("\n\nPage fault: cr2 0x%X, eip 0x%X, nos-err 0x%X\n", x86GetCr2(), r->eip, type);
50
51         HandleVirtFault(x86GetCr2(), type);
52     } else if (num == ISR_NMI) {
53         Panic(PANIC_NON_MASKABLE_INTERRUPT);
54     } else if (num == ISR_SYSTEM_CALL) {
55         r->eax = HandleSystemCall(r->eax, r->ebx, r->ecx, r->edx, r->esi, r->edi);
56     } else {
57         LogWriteSerial("Got interrupt %d. (r->eip = 0x%X)\n", num, r->eip);
58         UnhandledFault();
59     }
60 }
61
62 void ArchSendEoi(int irq_num) {
63     SendPicEoi(irq_num);
64 }
65
66 void ArchSetIrql(int irql) {
67     if (irql == IRQL_HIGH || irql == IRQL_TIMER || !x86IsReadyForIrqs()) {
68         /*
69          * Interrupts stay off.
70          */
71         return;
72     }
73
74     if (irql >= IRQL_DRIVER) {
75         int irq_num = irql - IRQL_DRIVER;
76
77         /*
78          * We want to disable all higher IRQs (as the PIC puts the lowest priority interrupts at
79          * high numbers), as well as our self. Allow IRQ2 to stay enabled as it is used internally.
80          */
81         uint16_t mask = (0xFFFF ^ ((1 << irq_num) - 1)) & ~(1 << 2);
82         DisablePicLines(mask);
83     } else {
84         /*
85          * Allow everything to go through.
86          */
87         DisablePicLines(0x0000);
88     }
89
90     ArchEnableInterrupts();
91 }
92
93 bool x86IsReadyForIrqs(void) {
94     return ready_for_irqs;
95 }
96
97 void x86MakeReadyForIrqs(void) {
98     ready_for_irqs = true;
99     RaiseIrql(GetIrql());
100 }
101
102
103
104
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  /*
6   * x86/dev/pic.c - Programmable Interrupt Controller
7   *
8   * The PIC controls the hardware raised interrupts (IRQs), i.e. those from
9   * 'external' devices such as the keyboard, system timer, disk, etc. Hence the
10  * PIC must be configured before any external interrupts can be seen by the CPU.
11  *
12  * There are two quirks to certain IRQs. The first comes about by the fact that
13  * each system actually has two PICs - the primary PIC handling IRQs 0-7, and the
14  * secondary PIC handling IRQs 8-15. They are connected (cascaded) to each other,
15  * using IRQ 2 for communication. Thus, an IRQ 2 will get to the CPU.
16  *
17  * The other quirk is the 'spurious interrupt', which can occur on IRQ 7 or 15.
18  * See pic_is_spurious for more details.
19  */
20
21  #define PIC1_COMMAND    0x20
22  #define PIC1_DATA       0x21
23  #define PIC2_COMMAND    0xA0
24  #define PIC2_DATA       0xA1
25
26  #define PIC_EOI         0x20
27  #define PIC_REG_ISR     0x0B
28
29  #define ICW1_ICW4       0x01
30  #define ICW1_INIT       0x10
31  #define ICW4_8086       0x01
32
33
34  /*
35   * Delay for a short period of time, for use in between IO calls to the PIC.
36   * This is required as some PICs have a hard time keeping up with the speed
37   * of modern CPUs (the original PIC was introduced in 1976!).
38   */
39  static void IoWait(void) {
40  }
41
42
43  /*
44   * Read an internal PIC register.
45   */
46  static uint16_t ReadPicReg(int ocw3) {
47      outb(PIC1_COMMAND, ocw3);
48      outb(PIC2_COMMAND, ocw3);
49      return ((uint16_t) inb(PIC2_COMMAND) << 8) | inb(PIC1_COMMAND);
50  }
51
52  /*
53   * Due to a race condition between the PIC and the CPU, we sometimes get a
54   * 'spurious' interrupt sent to the CPU on IRQ 7 or 15. If an IRQ 7 or 15
55   * arrives, we need to check if it is an actual interrupt or a spurious interrupt.
56   * Distinguishing them is important - spurious IRQs may cause drivers to misbehave,
57   * and we don't need to send an EOI after a spurious interrupt.
58   */
59  bool IsPicIrqSpurious(int irq_num) {
60      if (irq_num == PIC_IRQ_BASE + 7) {
61          uint16_t isr = ReadPicReg(PIC_REG_ISR);
62          if (!(isr & (1 << 7))) {
63              return true;
64          }
65      } else if (irq_num == PIC_IRQ_BASE + 15) {
66          uint16_t isr = ReadPicReg(PIC_REG_ISR);
67          if (!(isr & (1 << 15))) {
68              /*
69               * It is spurious, but the primary PIC doesn't know that, as it came
70               * from the secondary PIC. So only send an EOI to the primary PIC.
71               */
72              outb(PIC1_COMMAND, PIC_EOI);
73              return true;
74          }
75      }
76  }
77
78  return false;
79  }
80
81  /*
82   * Acknowledge the previous interrupt. We will not receive any interrupts of
83   * the same type until we have acknowledged it.
84   */
85  void SendPicEoi(int irq_num) {
86      if (irq_num >= PIC_IRQ_BASE + 8) {
87          outb(PIC2_COMMAND, PIC_EOI);
88      }
89
90      outb(PIC1_COMMAND, PIC_EOI);
91  }
92
93  /*
94   * Change which interrupt numbers are used by the IRQs. They will initially
95   * use interrupts 0 through 15, which isn't very good as it conflicts with the
96   * interrupt numbers for the CPU exceptions.
97   */
98  static void RemapPic(int offset) {
99      uint8_t mask1 = inb(PIC1_DATA);
100     uint8_t mask2 = inb(PIC2_DATA);
101
102     outb(PIC1_COMMAND, ICW1_INIT | ICW1_ICW4);
103     IoWait();
104     outb(PIC2_COMMAND, ICW1_INIT | ICW1_ICW4);
105     IoWait();
106     outb(PIC1_DATA, offset);
107     IoWait();
108     outb(PIC2_DATA, offset + 8);
109     IoWait();
110     outb(PIC1_DATA, 4);
111     IoWait();
112     outb(PIC2_DATA, 2);
113     IoWait();
114
115     outb(PIC1_DATA, ICW4_8086);
116     IoWait();
117     outb(PIC2_DATA, ICW4_8086);
118     IoWait();

```

```

119
120     outb(PIC1_DATA, mask1);
121     outb(PIC2_DATA, mask2);
122 }
123
124 /**
125  * Set which IRQ numbers are disabled. Overwrites the previous call to this function completely
126  * (i.e. this is an 'equals' operation, not an 'and' or 'or'.)
127  *
128  * @param irq_bitfield A bitfield of IRQs, where the lowest bit corresponds to IRQ0. For each bit,
129  *                     a zero will enable the interrupt, and a one will disable the interrupt.
130  *                     To disable all lines, specify 0xFFFF. To enable all lines, specify 0x0000.
131  */
132 void DisablePicLines(uint16_t irq_bitfield) {
133     static uint16_t prev = 0xFFFF;    // we initially set it to 0, so this will be different
134
135     if (prev == irq_bitfield) {
136         return;
137     }
138
139     outb(PIC1_DATA, irq_bitfield & 0xFF);
140     outb(PIC2_DATA, irq_bitfield >> 8);
141     prev = irq_bitfield;
142 }
143
144 /**
145  * Initialise the PIC.
146  */
147 void InitPic(void) {
148     /*
149      * Remap the PIC so that interrupts 0-15 are mapped to 32-47.
150      * This way we don't have conflicts with the CPU exceptions which are
151      * hard-wired to these interrupt numbers.
152      */
153     RemapPic(PIC_IRQ_BASE);
154
155     /*
156      * Now we can disable all masks, allowing interrupts to reach the CPU.
157      */
158     DisablePicLines(0x0000);
159 }

```

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     int res = RequireDriver(filename);
11     if (res != 0) {
12         PanicEx(PANIC_REQUIRED_DRIVER_NOT_FOUND, filename);
13     }
14
15     size_t addr = GetSymbolAddress(init);
16     if (addr == 0) {
17         PanicEx(PANIC_REQUIRED_DRIVER_MISSING_SYMBOL, filename);
18     }
19
20     return addr;
21 }
22
23 static void LoadSlowDriversInBackground(void*) {
24     /*
25      * To make the OS boot faster, we'll load the less critical, and slower
26      * drivers in a new thread. This means we can continue initialising the rest
27      * of the OS while drivers load.
28      */
29     ((void(*)()) (Loadx86Driver("sys/acpi.sys", "InitAcpica"))())();
30 }
31
32 void ArchInitDev(bool fs) {
33     if (!fs) {
34         InitIde();
35         //InitFloppy();
36     } else {
37         ((void(*)()) (Loadx86Driver("sys/vesa.sys", "InitVesa"))())();
38         ((void(*)()) (Loadx86Driver("sys/ps2.sys", "InitPs2"))())();
39         CreateThread(LoadSlowDriversInBackground, NULL, GetVas(), "drvloader");
40     }
41 }
42 }

```

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 <irq1.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* flp, int cmd) {
48     int base = flp->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     LogWriteSerial("floppy_write_cmd: timeout\n");
58 }
59
60 static uint8_t FloppyReadData(struct floppy_data* flp) {
61     int base = flp->base;
62     for (int i = 0; i < 60; ++i) {
63         SleepMilli(10);
64         if (inb(base + FLOPPY_MSR) & 0x80) {
65             return inb(base + FLOPPY_FIFO);
66         }
67     }
68     LogWriteSerial("floppy_read_data: timeout\n");
69     return 0;
70 }
71
72 static void FloppyCheckInterrupt(struct floppy_data* flp, int* st0, int* cyl) {
73     LogWriteSerial("FloppyCheckInterrupt\n");
74     FloppyWriteCommand(flp, CMD_SENSE_INT);
75     *st0 = FloppyReadData(flp);
76     *cyl = FloppyReadData(flp);
77 }
78
79 /*
80  * The state can be 0 (off), 1 (on) or 2 (currently on, but will shortly be turned off).
81  */
82
83 static volatile int floppy_motor_state = 0;
84 static volatile int floppy_motor_ticks = 0;
85
86 static void FloppyMotor(struct floppy_data* flp, bool state) {
87     LogWriteSerial("FloppyMotor\n");
88     int base = flp->base;
89
90     if (state) {
91         if (!floppy_motor_state) {
92             outb(base + FLOPPY_DOR, 0x1C);
93             SleepMilli(150);
94         }
95         floppy_motor_state = 1;
96     } else {
97         floppy_motor_state = 2;
98         floppy_motor_ticks = 1000;
99     }
100 }
101
102 static volatile bool floppy_got_irq = false;
103
104 static void FloppyIrqWait() {
105     LogWriteSerial("FloppyIrqWait\n");
106
107     /*
108      * Wait for the interrupt to come. If it doesn't the system will probably
109      * lockup.
110      */
111     while (!floppy_got_irq) {
112         SleepMilli(10);
113     }
114
115     /*
116      * Clear it for next time.
117      */
118     floppy_got_irq = false;
119 }
120
121 static int FloppyIrqHandler(struct x86_regs*) {
122     LogWriteSerial("FloppyIrqHandler\n");
123     floppy_got_irq = true;
124     return 0;
125 }
126
127 static void FloppyMotorControlThread(void*) {
128     LogWriteSerial("FloppyMotorControlThread\n");
129
130     while (1) {
131         SleepMilli(50);
132         if (floppy_motor_state == 2) {
133             floppy_motor_ticks -= 50;
134             if (floppy_motor_ticks <= 0) {
135                 /*
136                  * Actually turn off the motor.
137                  */
138                 //outb(0x3F0 + FLOPPY_DOR, 0x0C);
139                 //floppy_motor_state = 0;
140             }
141         }
142     }
143 }
144
145 /*
146  * Move to cylinder 0.
147  */
148 static int FloppyCalibrate(struct floppy_data* flp) {
149     LogWriteSerial("FloppyCalibrate\n");
150
151     int st0 = -1;
152     int cyl = -1;
153
154     FloppyMotor(flp, true);
155
156     for (int i = 0; i < 10; ++i) {
157         FloppyWriteCommand(flp, CMD_RECALIBRATE);
158         FloppyWriteCommand(flp, 0);
159     }
160 }

```

```

166     FloppyIrqWait();
167     FloppyCheckInterrupt(flp, &st0, &cyl);
168
169     if (st0 & 0xC0) {
170         continue;
171     }
172
173     if (cyl == 0) {
174         FloppyMotor(flp, false);
175         return 0;
176     }
177 }
178
179 LogDeveloperWarning("couldn't calibrate floppy\n");
180 FloppyMotor(flp, false);
181 return EIO;
182 }
183
184 static void FloppyConfigure(struct floppy_data* flp) {
185     LogWriteSerial("FloppyConfigure\n");
186
187     FloppyWriteCommand(flp, CMD_CONFIGURE);
188     FloppyWriteCommand(flp, 0x00);
189     FloppyWriteCommand(flp, 0x08);
190     FloppyWriteCommand(flp, 0x00);
191 }
192
193 static int FloppyReset(struct floppy_data* flp) {
194     LogWriteSerial("FloppyReset\n");
195
196     int base = flp->base;
197     outb(base + FLOPPY_DOR, 0x00);
198     outb(base + FLOPPY_DOR, 0x0C);
199
200     FloppyIrqWait();
201
202     for (int i = 0; i < 4; ++i) {
203         int st0, cyl;
204         FloppyCheckInterrupt(flp, &st0, &cyl);
205     }
206
207     outb(base + FLOPPY_CCR, 0x00);
208
209     FloppyMotor(flp, true);
210
211     FloppyWriteCommand(flp, CMD_SPECIFY);
212     FloppyWriteCommand(flp, 0xDF);
213     FloppyWriteCommand(flp, 0x02);
214
215     SleepMilli(300);
216     FloppyConfigure(flp);
217     SleepMilli(300);
218     FloppyMotor(flp, false);
219
220     int res = FloppyCalibrate(flp);
221     LogWriteSerial("FloppyReset: res = %d\n", res);
222     return res;
223 }
224
225 static int FloppySeek(struct floppy_data* flp, int cylinder, int head) {
226     LogWriteSerial("FloppySeek\n");
227
228     int st0, cyl;
229     FloppyMotor(flp, true);
230
231     for (int i = 0; i < 10; ++i) {
232         FloppyWriteCommand(flp, CMD_SEEK);
233         FloppyWriteCommand(flp, head < 2);
234         FloppyWriteCommand(flp, cylinder);
235
236         FloppyIrqWait();
237         FloppyCheckInterrupt(flp, &st0, &cyl);
238
239         if (st0 & 0xC0) {
240             continue;
241         }
242
243         if (cyl == cylinder) {
244             FloppyMotor(flp, false);
245             return 0;
246         }
247     }
248
249     LogWriteSerial("couldn't seek floppy\n");
250     FloppyMotor(flp, false);
251     return EIO;
252 }
253
254 static void FloppyDmaInit(void) {
255     LogWriteSerial("FloppyDmaInit\n");
256
257     /*
258     * Put the data at *physical address* 0x10000. The address can be anywhere
259     * under 24MB that doesn't cross a 64KB boundary. We choose this location as
260     * it should be unused as this is where the temporary copy of the kernel was
261     * stored during boot.
262     */
263     uint32_t addr = (uint32_t) 0x10000;
264
265     /*
266     * We must give the DMA the actual count minus 1.
267     */
268     int count = 0x4800 - 1;
269
270     /*
271     * Send some magical stuff to the DMA controller.
272     */
273     outb(0x0A, 0x06);
274     outb(0x0C, 0xFF);
275     outb(0x04, (addr >> 0) & 0xFF);
276     outb(0x04, (addr >> 8) & 0xFF);
277     outb(0x81, (addr >> 16) & 0xFF);
278     outb(0x0C, 0xFF);
279     outb(0x05, (count >> 0) & 0xFF);
280     outb(0x05, (count >> 8) & 0xFF);
281     outb(0x0B, 0x46);
282     outb(0x0A, 0x02);
283
284
285 static int FloppyDoCylinder(struct floppy_data* flp, int cylinder) {
286     LogWriteSerial("FloppyDoCylinder\n");
287
288     /*
289     * Move both heads to the correct cylinder.
290     */
291     if (FloppySeek(flp, cylinder, 0) != 0) return EIO;
292     if (FloppySeek(flp, cylinder, 1) != 0) return EIO;
293
294     /*
295     * This time, we'll try up to 20 times.

```

```

296 */
297 for (int i = 0; i < 20; ++i) {
298     LogWriteSerial("READ ATTEMPT %d\n", i + 1);
299     FloppyMotor(flp, true);
300
301     if (i % 5 == 3) {
302         if (i % 10 == 8) {
303             LogWriteSerial("resetting floppy!\n");
304             FloppyReset(flp);
305             FloppyMotor(flp, true);
306         }
307
308         FloppyCalibrate(flp);
309
310         if (FloppySeek(flp, cylinder, 0) != 0) return EIO;
311         if (FloppySeek(flp, cylinder, 1) != 0) return EIO;
312     }
313
314     FloppyDmaInit();
315
316     SleepMilli(100);
317
318     /*
319     * Send the read command.
320     */
321     FloppyWriteCommand(flp, CMD_READ | 0xC0);
322     FloppyWriteCommand(flp, 0);
323     FloppyWriteCommand(flp, cylinder);
324     FloppyWriteCommand(flp, 0);
325     FloppyWriteCommand(flp, 1);
326     FloppyWriteCommand(flp, 2);
327     FloppyWriteCommand(flp, 18);
328     FloppyWriteCommand(flp, 0x1B);
329     FloppyWriteCommand(flp, 0xFF);
330
331     FloppyIrqWait();
332
333     /*
334     * Read back some status information, some of which is very mysterious.
335     */
336     uint8_t st0, st1, st2, rcy, rhe, rse, bps;
337     st0 = FloppyReadData(flp);
338     st1 = FloppyReadData(flp);
339     st2 = FloppyReadData(flp);
340     rcy = FloppyReadData(flp);
341     rhe = FloppyReadData(flp);
342     rse = FloppyReadData(flp);
343     bps = FloppyReadData(flp);
344
345     /*
346     * Check for errors. More tests can be done, but it would make the code
347     * even longer.
348     */
349     if (st0 & 0xC0) {
350         static const char * status[] = { 0, "error", "invalid command", "drive not ready" };
351         LogWriteSerial("floppy_do_sector: status = %s\n", status[st0 >> 6]);
352         LogWriteSerial("st0 = 0x%X, st1 = 0x%X, st2 = 0x%X, bps = %d\n", st0, st1, st2, bps);
353         continue;
354     }
355     if (st1 & 0x80) {
356         continue;
357     }
358     if (st0 & 0x8) {
359         continue;
360     }
361     if (st1 & 0x20) {
362         continue;
363     }
364     if (st1 & 0x10) {
365         continue;
366     }
367     if (bps != 2) {
368         continue;
369     }
370
371     (void) st2;
372     (void) rcy;
373     (void) rhe;
374     (void) rse;
375
376     FloppyMotor(flp, false);
377     return 0;
378 }
379
380 LogWriteSerial("couldn't read floppy\n");
381 FloppyMotor(flp, false);
382 return EIO;
383 }
384
385
386
387 static int FloppyIo(struct floppy_data* flp, struct transfer* io) {
388     LogWriteSerial("FloppyIo\n");
389     EXACT_IOCTL(IRQL_STANDARD);
390
391     if (io->direction == TRANSFER_WRITE) {
392         return EROFS;
393     }
394
395     int lba = io->offset / 512;
396     int count = io->length_remaining / 512;
397
398     if (io->offset % 512 != 0) {
399         return EINVAL;
400     }
401     if (io->length_remaining % 512 != 0) {
402         return EINVAL;
403     }
404     if (count <= 0 || count > 0xFF || lba < 0 || lba >= 2880) {
405         return EINVAL;
406     }
407
408     AcquireMutex(floppy_lock, -1);
409
410 next_sector:;
411     /*
412     * Floppies use CHS (cylinder, head, sector) for addressing sectors instead of
413     * LBA (linear block addressing). Hence we need to convert to CHS.
414     *
415     * Head: which side of the disk it is on (i.e. either the top or bottom)
416     * Cylinder: which 'slice' (cylinder) of the disk we should look at
417     * Sector: which 'ring' (sector) of that cylinder we should look at
418     *
419     * Note that sector is 1-based, whereas cylinder and head are 0-based.
420     * Don't ask why.
421     */
422     int head = (lba % (18 * 2)) / 18;
423     int cylinder = (lba / (18 * 2));
424     int sector = (lba % 18) + 1;
425

```

```

426  /*
427  * Cylinder 0 has some commonly used data, so cache it seperately for improved
428  * speed.
429  */
430  if (cylinder == 0) {
431      if (!flp->got_cylinder_zero) {
432          flp->got_cylinder_zero = true;
433
434          int error = FloppyDoCylinder(flp, cylinder);
435
436          if (error != 0) {
437              ReleaseMutex(floppy_lock);
438              return error;
439          }
440
441          memcpy(flp->cylinder_zero, flp->cylinder_buffer, 0x4800);
442
443          /*
444           * Must do this as we trashed the cached cylinder.
445           * Luckily we only ever need to do this once.
446           */
447          flp->stored_cylinder = cylinder;
448      }
449
450      PerformTransfer(flp->cylinder_zero + (512 * (sector - 1 + head * 18)), io, 512);
451  }
452  else {
453      if (cylinder != flp->stored_cylinder) {
454          int error = FloppyDoCylinder(flp, cylinder);
455
456          if (error != 0) {
457              ReleaseMutex(floppy_lock);
458              return error;
459          }
460      }
461
462      PerformTransfer(flp->cylinder_buffer + (512 * (sector - 1 + head * 18)), io, 512);
463      flp->stored_cylinder = cylinder;
464  }
465
466  /*
467  * Read only read one sector, so we need to repeat the process if multiple sectors
468  * were requested. We could just do a larger copy above, but this is a bit simpler,
469  * as we don't need to worry about whether the entire request is on cylinder or not.
470  */
471  if (count > 0) {
472      ++lba;
473      --count;
474      goto next_sector;
475  }
476
477  ReleaseMutex(floppy_lock);
478  return 0;
479
480 static bool IsSeekable(struct vnode*) {
481     return true;
482 }
483
484 static int ReadWrite(struct vnode* node, struct transfer* io) {
485     return FloppyIo(node->data, io);
486 }
487
488 static int Create(struct vnode* node, struct vnode** partition, const char* name, int, mode_t) {
489     AcquireMutex(floppy_lock, -1);
490     struct floppy_data* flp = node->data;
491     int res = DiskCreateHelper(&flp->partitions, partition, name);
492     ReleaseMutex(floppy_lock);
493     return res;
494 }
495
496 static int Follow(struct vnode* node, struct vnode** output, const char* name) {
497     AcquireMutex(floppy_lock, -1);
498     struct floppy_data* flp = node->data;
499     int res = DiskFollowHelper(&flp->partitions, output, name);
500     ReleaseMutex(floppy_lock);
501     return res;
502 }
503
504 static uint8_t DirentType(struct vnode*) {
505     return DT_BLK;
506 }
507
508 static int Stat(struct vnode*, struct stat* st) {
509     st->st_mode = S_IFBLK | S_IRWXU | S_IRWXG | S_IRWXO;
510     st->st_atime = 0;
511     st->st_blksize = 512;
512     st->st_blocks = 2880;
513     st->st_ctime = 0;
514     st->st_dev = 0xBA9E9E9E;
515     st->st_gid = 0;
516     st->st_ino = 0xCAFEBABE;
517     st->st_mtime = 0;
518     st->st_nlink = 1;
519     st->st_rdev = 0xCAFEBABE;
520     st->st_size = 1024 * 1440;
521     st->st_uid = 0;
522     return 0;
523 }
524
525 static int Close(struct vnode*) {
526     return 0;
527 }
528
529 static const struct vnode_operations dev_ops = {
530     .is_seekable = IsSeekable,
531     .read = ReadWrite,
532     .write = ReadWrite,
533     .close = Close,
534     .create = Create,
535     .follow = Follow,
536     .dirent_type = DirentType,
537     .stat = Stat,
538 };
539
540 void InitFloppy(void) {
541     floppy_lock = CreateMutex("floppy");
542     CreateThread(FloppyMotorControlThread, NULL, GetVas(), "flpmotor");
543
544     for (int i = 0; i < 1; ++i) {
545         struct vnode* node = CreateVnode(dev_ops);
546         struct floppy_data* flp = AllocHeap(sizeof(struct floppy_data));
547
548         RegisterIrqHandler(PIC_IRQ_BASE + 6, FloppyIrqHandler);
549
550         flp->disk_num = 0;
551         flp->base = 0x3F0;
552         flp->cylinder_buffer = (uint8_t*) MapVirt(0, 0, CYLINDER_SIZE, VM_READ | VM_WRITE | VM_LOCK, NULL, 0);
553         flp->cylinder_zero = (uint8_t*) MapVirt(0, 0, CYLINDER_SIZE, VM_READ | VM_WRITE | VM_LOCK, NULL, 0);
554     }
555 }

```



```

556     flp->got_cylinder_zero = false;
557     flp->stored_cylinder   = -1;
558     FloppyReset(flp);
559
560     InitDiskPartitionHelper(&flp->partitions);
561
562     node->data = flp;
563     AddVfsMount(node, GenerateNewRawDiskName(DISKUTIL_TYPE_FLOPPY));
564     CreateDiskPartitions(CreateOpenFile(node, 0, 0, true, true));
565 }
566

```

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 <irq.h>
16 #include <diskutil.h>
17
18 #define MAX_TRANSFER_SIZE (1024 * 16)
19
20 struct semaphore* ide_lock = NULL;
21
22 struct ide_data {
23     int disk_num;
24     unsigned int sector_size;
25     uint64_t total_num_sectors;
26     uint16_t* transfer_buffer;
27     size_t primary_base;
28     size_t primary_alternative;
29     size_t secondary_base;
30     size_t secondary_alternative;
31     size_t busmaster_base;
32
33     struct disk_partition_helper partitions;
34 };
35
36 int IdeCheckError(struct ide_data* ide) {
37     uint16_t base = ide->disk_num >= 2 ? ide->secondary_base : ide->primary_base;
38
39     uint8_t status = inb(base + 0x7);
40     if (status & 0x01) {
41         return EIO;
42     } else if (status & 0x20) {
43         return EIO;
44     } else if (!(status & 0x08)) {
45         return EIO;
46     }
47     return 0;
48 }
49
50 int IdePoll(struct ide_data* ide) {
51     uint16_t base = ide->disk_num >= 2 ? ide->secondary_base : ide->primary_base;
52     uint16_t alt_status_reg = ide->disk_num >= 2 ? ide->secondary_alternative : ide->primary_alternative;
53
54     /*
55      * Delay for a moment by reading the alternate status register.
56      */
57     for (int i = 0; i < 4; ++i) {
58         inb(alt_status_reg);
59     }
60
61     /*
62      * Wait for the device to not be busy. We have a timeout in case the
63      * device is faulty, we don't want to be in an endless loop and freeze
64      * the kernel.
65      */
66     int timeout = 0;
67     while (inb(base + 0x7) & 0x80) {
68         if (timeout > 975) {
69             SleepMilli(10);
70         }
71         if (timeout++ > 1000) {
72             return EIO;
73         }
74     }
75     return 0;
76 }
77
78 /*
79  * Read or write the primary ATA drive on the first controller. We use LBA28,
80  * so we are limited to a 28 bit sector number (i.e. disks up to 128GB in size)
81  */
82 static int IdeIo(struct ide_data* ide, struct transfer* io) {
83     EXACT_IRQL(IRQL_STANDARD);
84     int disk_num = ide->disk_num;
85
86     /*
87      * IDE devices do not contain an (accessible) disk buffer in PIO mode, as
88      * they transfer data through the IO ports. Hence we must read/write into
89      * this buffer first, and then move it safely to the destination.
90      * (we could use DMA instead, but PIO is simpler)
91      */
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     if (count <= 0 || sector < 0 || sector > 0xFFFFFFF || (uint64_t) sector + count >= ide->total_num_sectors) {
108         return EINVAL;
109     }
110 }

```

```

113
114 AcquireSemaphore(&ide_lock, -1);
115
116 uint16_t base = disk_num >= 2 ? ide->secondary_base : ide->primary_base;
117 uint16_t dev_ctrl_reg = disk_num >= 2 ? ide->secondary_alternative : ide->primary_alternative;
118
119 int max_sectors_at_once = MAX_TRANSFER_SIZE / ide->sector_size;
120 if (max_sectors_at_once > 255) {
121     // Hardware limitation
122     max_sectors_at_once = 255;
123 }
124
125 while (count > 0) {
126     int sectors_in_this_transfer = count > max_sectors_at_once ? max_sectors_at_once : count;
127
128     if (io->direction == TRANSFER_WRITE) {
129         PerformTransfer(buffer, io, ide->sector_size);
130     }
131
132     /*
133     * Send a whole heap of flags and the high 4 bits of the LBA to the controller.
134     */
135     outb(base + 0x6, 0xE0 | ((disk_num & 1) << 4) | ((sector >> 24) & 0xF));
136
137     /*
138     * Disable interrupts, we are going to use polling.
139     */
140     outb(dev_ctrl_reg, 2);
141
142     /*
143     * May not be needed, but it doesn't hurt to do it.
144     */
145     outb(base + 0x1, 0x00);
146
147     /*
148     * Send the number of sectors, and the sector's LBA.
149     */
150     outb(base + 0x2, sectors_in_this_transfer);
151     outb(base + 0x3, (sector >> 0) & 0xFF);
152     outb(base + 0x4, (sector >> 8) & 0xFF);
153     outb(base + 0x5, (sector >> 16) & 0xFF);
154
155     /*
156     * Send either the read or write command.
157     */
158     outb(base + 0x7, io->direction == TRANSFER_WRITE ? 0x30 : 0x20);
159
160     /*
161     * Wait for the data to be ready.
162     */
163     IdePoll(ide);
164
165     /*
166     * Read/write the data from/to the disk using ports.
167     */
168     if (io->direction == TRANSFER_WRITE) {
169         for (int c = 0; c < sectors_in_this_transfer; ++c) {
170             if (c != 0) {
171                 PerformTransfer(buffer, io, ide->sector_size);
172                 IdePoll(ide);
173             }
174
175             for (uint64_t i = 0; i < ide->sector_size / 2; ++i) {
176                 outw(base + 0x00, buffer[i]);
177             }
178             IdePoll(ide);
179
180             /*
181             * We need to flush the disk's cache if we are writing.
182             */
183             outb(base + 0x7, 0xE7);
184             IdePoll(ide);
185
186         } else {
187             int err = IdeCheckError(ide);
188             if (err) {
189                 ReleaseSemaphore(&ide_lock);
190                 return err;
191             }
192
193             for (int c = 0; c < sectors_in_this_transfer; ++c) {
194                 if (c != 0) {
195                     IdePoll(ide);
196                 }
197
198                 for (uint64_t i = 0; i < ide->sector_size / 2; ++i) {
199                     buffer[i] = inw(base + 0x00);
200                 }
201                 PerformTransfer(buffer, io, ide->sector_size);
202             }
203         }
204
205         /*
206         * Get ready for the next part of the transfer.
207         */
208         count -= sectors_in_this_transfer;
209         sector += sectors_in_this_transfer;
210     }
211 }
212
213 ReleaseSemaphore(&ide_lock);
214
215 return 0;
216
217 static int IdeGetNumSectors(struct ide_data* ide) {
218     AcquireSemaphore(&ide_lock, -1);
219
220     uint16_t base = ide->disk_num >= 2 ? ide->secondary_base : ide->primary_base;
221
222     /*
223     * Select the correct drive.
224     */
225     outb(base + 0x6, 0xE0 | ((ide->disk_num & 1) << 4));
226
227     /*
228     * Send the READ NATIVE MAX ADDRESS command, which will return the size
229     * of disk in sectors.
230     */
231     outb(base + 0x7, 0xF8);
232
233     IdePoll(ide);
234
235     /*
236     * The outputs are in the same registers we use to put the LBA
237     * when we read/write from the disk.
238     */
239     int sectors = 0;
240     sectors |= (int) inb(base + 0x3);
241     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 bool IsSeekable(struct vnode*) {
252     return true;
253 }
254
255 static int ReadWrite(struct vnode* node, struct transfer* io) {
256     return IdeIo(node->data, io);
257 }
258
259 static int Create(struct vnode* node, struct vnode** partition, const char* name, int, mode_t) {
260     AcquireSemaphore(ide_lock, -1);
261     struct ide_data* ide = node->data;
262     int res = DiskCreateHelper(ide->partitions, partition, name);
263     ReleaseSemaphore(ide_lock);
264     return res;
265 }
266
267 static int Follow(struct vnode* node, struct vnode** output, const char* name) {
268     AcquireSemaphore(ide_lock, -1);
269     struct ide_data* ide = node->data;
270     int res = DiskFollowHelper(ide->partitions, output, name);
271     ReleaseSemaphore(ide_lock);
272     return res;
273 }
274
275 static uint8_t DirentType(struct vnode*) {
276     return DT_BLK;
277 }
278
279 static int Stat(struct vnode* node, struct stat* st) {
280     struct ide_data* ide = node->data;
281
282     st->st_mode = S_IFBLK | S_IRWXU | S_IRWXG | S_IRWXO;
283     st->st_atime = 0;
284     st->st_blksize = ide->sector_size;
285     st->st_blocks = ide->total_num_sectors;
286     st->st_ctime = 0;
287     st->st_dev = 0xBAFECAFE;
288     st->st_gid = 0;
289     st->st_ino = 0xCAFEBAFE;
290     st->st_mtime = 0;
291     st->st_nlink = 1;
292     st->st_rdev = 0xCAFEDEAD;
293     st->st_size = ide->sector_size * ide->total_num_sectors;
294     st->st_uid = 0;
295
296     return 0;
297 }
298
299 static int Close(struct vnode*) {
300     return 0;
301 }
302
303 static const struct vnode_operations dev_ops = {
304     .is_seekable = IsSeekable,
305     .read = ReadWrite,
306     .write = ReadWrite,
307     .close = Close,
308     .create = Create,
309     .follow = Follow,
310     .dirent_type = DirentType,
311     .stat = Stat,
312 };
313
314 void InitIde(void) {
315     ide_lock = CreateMutex("ide");
316
317     for (int i = 0; i < 1; ++i) {
318         struct vnode* node = CreateVnode(dev_ops);
319         struct ide_data* ide = AllocHeap(sizeof(struct ide_data));
320
321         ide->disk_num = 0;
322         ide->primary_base = 0x1F0;
323         ide->secondary_base = 0x170;
324         ide->primary_alternative = 0x3F6;
325         ide->secondary_alternative = 0x376;
326         ide->busmaster_base = 0x0;
327         ide->sector_size = 512;
328         ide->total_num_sectors = IdeGetNumSectors(ide);
329         ide->transfer_buffer = (uint16_t*) MapVirt(0, 0, MAX_TRANSFER_SIZE, VM_READ | VM_WRITE | VM_LOCK, NULL, 0);
330         InitDiskPartitionHelper(ide->partitions);
331
332         node->data = ide;
333         AddVfsMount(node, GenerateNewRawDiskName(DISKUTIL_TYPE_FIXED));
334         CreateDiskPartitions(CreateOpenFile(node, 0, 0, true, true));
335     }
336 }

```

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 /*
13 * x86/dev/pit.c - Programmable Interval Timer
14 *
15 * A system timer which can generate an interrupt (on IRQ 0) at regular intervals.
16 */
17
18 static uint64_t pit_hertz = 0;
19 static uint64_t pit_nanos = 0;
20
21 static int HandlePit(struct x86_regs* r) {
22     (void) r;
23     assert(pit_nanos != 0);
24
25     ReceivedTimer(pit_nanos);
26
27     return 0;
28 }
29
30 void InitPit(int hertz) {
31     pit_hertz = (uint64_t) hertz;
32
33     int divisor = 1193180 / hertz;
34     outb(0x43, 0x36);
35     outb(0x40, divisor & 0xFF);
36     outb(0x40, divisor >> 8);
37
38     pit_nanos = 1000000000ULL / pit_hertz;
39
40     RegisterIrqHandler(PIC_IRQ_BASE + 0, HandlePit);
41 }

```

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/idt.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     lidt [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 irqs[])() = {
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) irqs[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     MAX_IRQL(IRQL_PAGE_FAULT);
314
315     off_t file_size;
316     int res = GetFileSize(file, &file_size);
317     if (res != 0) {
318         return res;
319     }
320
321     size_t file_rgn = MapVirt(0, 0, file_size, VM_READ | VM_FILE, file, 0);
322     res = ElfLoad((void*) file_rgn, relocation_point, file, table);
323
324     struct Elf32_Ehdr* elf_header = (struct Elf32_Ehdr*) file_rgn;
325     struct Elf32_Shdr* sect_headers = (struct Elf32_Shdr*) (file_rgn + elf_header->e_shoff);
326
327     for (int i = 0; i < elf_header->e_shnum; ++i) {
328         const char* sh_name = ElfLookupString((void*) file_rgn, sect_headers[i].sh_name);
329         if (!strcmp(sh_name, ".lockedtext") || !strcmp(sh_name, ".lockeddata")) {
330             // it's okay to lock extra memory - it wouldn't be ok if we did it the other way around though
331             // (assumed everything was locked, and only unlocked parts of it)
332             size_t start_addr = (sect_headers[i].sh_addr - 0xD0000000U + *relocation_point) & ~(ARCH_PAGE_SIZE - 1);
333             size_t num_pages = (sect_headers[i].sh_size + ARCH_PAGE_SIZE - 1) / ARCH_PAGE_SIZE;
334             while (num_pages--)
335                 LockVirt(start_addr);
336             start_addr += ARCH_PAGE_SIZE;
337         }
338     }
339
340     UnmapVirt(file_rgn, file_size);
341
342     return res;
343 }
344
345 void ArchLoadSymbols(struct open_file* file, size_t adjust) {
346     off_t size;
347     int res = GetFileSize(file, &size);
348     if (res != 0)
349         Panic(PANIC_BAD_KERNEL);
350
351     size_t mem = MapVirt(0, 0, size, VM_READ | VM_FILE, file, 0);
352
353     struct Elf32_Ehdr* elf_header = (struct Elf32_Ehdr*) mem;
354
355     /*
356      * These should never happen - otherwise our kernel shouldn't be running!
357      */
358     if (!IsElfValid(elf_header) || elf_header->e_shoff == 0) {
359         Panic(PANIC_BAD_KERNEL);
360     }
361
362     struct Elf32_Shdr* section_headers = (struct Elf32_Shdr*) (size_t) (mem + elf_header->e_shoff);
363     size_t symbol_table_offset = 0;
364     size_t symbol_table_length = 0;
365     size_t string_table_offset = 0;
366     size_t string_table_length = 0;
367
368     /*
369      * Find the address and size of the symbol and string tables.
370      */
371     for (int i = 0; i < elf_header->e_shnum; ++i) {
372         size_t file_offset = (section_headers + i)->sh_offset;
373         size_t address = (section_headers + i)->sh_offset + (section_headers + i)->sh_name;
374
375         char* name_buffer = (char*) (mem + address);
376
377         if (!strcmp(name_buffer, ".symtab")) {
378             symbol_table_offset = file_offset;
379             symbol_table_length = (section_headers + i)->sh_size;
380         } else if (!strcmp(name_buffer, ".strtab")) {
381             string_table_offset = file_offset;
382             string_table_length = (section_headers + i)->sh_size;
383         }
384     }
385
386     if (symbol_table_offset == 0 || string_table_offset == 0 || symbol_table_length == 0 || string_table_length == 0) {
387         Panic(PANIC_BAD_KERNEL);
388     }
389
390     struct Elf32_Sym* symbol_table = (struct Elf32_Sym*) (mem + symbol_table_offset);
391     const char* string_table = (const char*) (mem + string_table_offset);
392
393     /*
394      * Register all of the visible symbols we find.
395      */
396     for (size_t i = 0; i < symbol_table_length / sizeof(struct Elf32_Sym); ++i) {
397         struct Elf32_Sym symbol = symbol_table[i];
398
399         if (symbol.st_value == 0) {
400             continue;
401         }
402
403         /*
404          * Skip "hidden" and "internal" symbols
405          */
406         if ((symbol.st_other & 3) != 0) {
407             continue;
408         }
409
410         /*
411          * No need for strdup, as the symbol table will call strdup anyway.
412          */
413         AddSymbol(string_table + symbol.st_name, symbol.st_value + adjust);
414     }
415
416     UnmapVirt(mem, size);
417 }
418
419
420

```

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[1];
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 ? 0xFFFFFFFFFU : 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 + 0xFFFF) & ~0xFFFF;
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 < 0x1000000) {
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;
138 }

```

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 <awi.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 / 0x4000000;
43     size_t page_num = (virtual % 0x4000000) / 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 / 0x4000000;
65     size_t page_num = (virtual % 0x4000000) / 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 + 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 |
160 kernel_page_directory[768] = ((size_t) first_page_table - 0xC0000000) | x86_PAGE_PRESENT | x86_PAGE_WRITE | x86_PAGE_USER;
161 |
162 /* <= is required to make it match kernel_entry.s */
163 size_t num_pages = (max_kernel_addr - 0xC0000000) / ARCH_PAGE_SIZE;
164 for (size_t i = 0; i < num_pages; ++i) {
165     first_page_table[i] = (i * ARCH_PAGE_SIZE) | x86_PAGE_PRESENT | x86_PAGE_WRITE;
166 }
167 |
168 /*
169  * Set up recursive mapping by mapping the 1024th page table to
170  * the page directory. See arch_vas_set_entry for an explanation of why we do this.
171  * "Locking" this page directory entry is the only way we can lock the final page of virtual
172  * memory, due to the recursive nature of this entry.
173  */
174 kernel_page_directory[1023] = ((size_t) kernel_page_directory - 0xC0000000) | x86_PAGE_PRESENT | x86_PAGE_WRITE;
175 |
176 vas->arch_data->p_page_directory = ((size_t) kernel_page_directory) - 0xC0000000;
177 vas->arch_data->v_page_directory = kernel_page_directory;
178 |
179 SetVas(vas);
180 |
181 /*
182  * The virtual memory manager is now initialised, so we can fill in
183  * the rest of the kernel state. This is important as we need all of
184  * the kernel address spaces to share page tables, so we must allocate
185  * them all now so new address spaces can copy from us.
186  */
187 /*for (int i = 769; i < 1023; ++i) {
188     x86AllocatePageTable(vas, i);
189 }*/
190 |
191 /*
192  * The maximum amount of virtual kernel memory we can access will depend on
193  * the amount of RAM - full access requires 1MB, which is an issue if we've got, e.g.
194  * only 1.5MB of RAM. We ensure we always get at least 128MB of kernel virtual memory -
195  * systems with 1.5MB of RAM will certainly not need that much virtual memory.
196  */
197 * A quick reference table:
198 *
199 * Physical RAM   Max Kernel Virtual RAM   Physical RAM Usage
200 * 1280 KB       132 MB                     248 / 544 (54% free)
201 * 1536 KB       194 MB                     312 / 800 (61% free)
202 * 2048 KB       324 MB                     440 / 1312 (66% free)
203 * 3072 KB       580 MB                     696 / 2336 (70% free)
204 * 4096 KB       764 MB                     880 / 3360 (73% free)
205 * 8192 KB       764 MB                     884 / 7456 (88% free)
206 */
207 |
208 size_t tables_allocated = 0;
209 int start = ARCH_KRNL_SBRK_BASE // 0x4000000;
210 for (int i = start; i < 1023; ++i) {
211     if (i >= start + 32 * tables_allocated * 16 > GetTotalPhysKilobytes()) {
212         break;
213     }
214     ++tables_allocated;
215     x86AllocatePageTable(vas, i);
216 }
217 |
218 LogWriteSerial("can access %d MB of kernel virtual memory\n", tables_allocated * 4);
219 |
220 /*
221  * The boot assembly code set up two page tables for us, that we no longer need.
222  * We can release that physical memory.
223  */
224 extern size_t boot_page_directory;
225 extern size_t boot_page_table1;
226 DeallocPhys(ArchVirtualToPhysical((size_t) &boot_page_directory));
227 DeallocPhys(ArchVirtualToPhysical((size_t) &boot_page_table1));
228 |

```

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 addr;
23     size_t size;
24     size_t lba;
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 void SynchroniseEntry(struct cache_entry* entry) {
42     // TODO: write to disk...
43     (void) entry;
44 }
45
46 static int Read(struct vnode*, struct transfer*) {
47     return 0;
48 }
49
50 static int Write(struct vnode*, struct transfer*) {
51     return 0;
52 }
53
54 static uint8_t DirentType(struct vnode*) {
55     return DT_BLK;
56 }
57
58 static int Stat(struct vnode* node, struct stat* st) {
59     struct cache_data* data = node->data;
60     return VnodeOpStat(data->underlying_disk->node, st);
61 }
62
63 static void TossCache(struct cache_data* data) {
64     AcquireMutex(data->lock, -1);
65     AvlTreeDestroy(data->cache);
66     data->cache = AvlTreeCreate();
67     ReleaseMutex(data->lock);
68 }
69
70 static void ReduceCache(struct cache_data* data) {
71     AcquireMutex(data->lock, -1);
72     // .. TODO: do something here...
73     ReleaseMutex(data->lock);
74 }
75
76 static int Close(struct vnode* node) {
77     TossCache(node->data);
78     return 0;
79 }
80
81 static const struct vnode_operations dev_ops = {
82     .read      = Read,
83     .write     = Write,
84     .dirent_type = DirentType,
85     .stat      = Stat,
86     .close     = Close
87 };
88
89 void RemoveCacheEntryHandler(void* entry_) {
90     struct cache_entry* entry = entry_;
91     SynchroniseEntry(entry);
92     UnmapVirt(entry->addr, entry->size);
93 }
94
95 struct open_file* CreateDiskCache(struct open_file* underlying_disk)
96 {
97     struct vnode* node = CreateVnode(dev_ops);
98     struct cache_data* data = AllocHeap(sizeof(struct cache_data));
99     data->underlying_disk = underlying_disk;
100     data->cache = AvlTreeCreate();
101     data->lock = CreateMutex("vcache");
102     data->block_size = 4096; // TODO: max of block size and ARCH_PAGE_SIZE
103     AvlTreeSetDeletionHandler(data->cache, RemoveCacheEntryHandler);
104     node->data = data;
105
106     struct open_file* cache = CreateOpenFile(node, underlying_disk->initial_mode, underlying_disk->flags, underlying_disk->can_read, underlying_disk->can_write);
107
108     AcquireMutex(cache_list_lock, -1);
109     LinkedListInsertEnd(cache_list, cache);
110     ReleaseMutex(cache_list_lock);
111
112     return cache;
113 }
114
115 static void ReduceCacheAmounts(bool toss) {
116     struct linked_list_node* node = LinkedListGetFirstNode(cache_list);
117     while (node != NULL) {
118         struct cache_data* data = ((struct open_file*) LinkedListGetDataFromNode(node))->node->data;
119         (toss ? TossCache : ReduceCache)(data);
120         node = LinkedListGetNextNode(node);
121     }
122 }
123
124 void SetDiskCaches(int mode) {
125     /*
126      * The PMM calls on allocation / free this before InitDiskCaches is called,
127      * so need to guard here.
128      */
129     if (cache_list_lock == NULL) {
130         return;
131     }
132
133     AcquireMutex(cache_list_lock, -1);
134
135     if (mode == DISKCACHE_REDUCE && current_mode == DISKCACHE_NORMAL) {
136         ReduceCacheAmounts(false);
137     } else if (mode == DISKCACHE_TOSS && current_mode != DISKCACHE_TOSS) {

```

```

139     ReduceCacheAmounts(true);
140 }
141
142 current_mode = mode;
143 ReleaseMutex(cache_list_lock);
144
145
146 void InitDiskCaches(void) {
147     cache_list = LinkedListCreate();
148     cache_list_lock = CreateMutex("vclist");
149 }

```

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 bool IsSeekable(struct vnode*) {
59     return true;
60 }
61
62 static int Create(struct vnode* node, struct vnode** fs, const char*, int flags, mode_t mode) {
63     struct partition_data* partition = node->data;
64     if (partition->fs != NULL) {
65         return EALREADY;
66     }
67
68     partition->fs = CreateOpenFile(*fs, flags, mode, true, true);
69     return 0;
70 }
71
72 static uint8_t DentryType(struct vnode*) {
73     return DT_BLK;
74 }
75
76 static int Stat(struct vnode* node, struct stat* st) {
77     struct partition_data* partition = node->data;
78
79     LogWriteSerial("calling stat on a partition... bps = %d, len = %d\n", partition->disk_bytes_per_sector, partition->length_bytes);
80
81     st->st_mode = S_IFBLK | S_IRWXU | S_IRWXG | S_IRWXO;
82     st->st_atime = 0;
83     st->st_blksize = partition->disk_bytes_per_sector;
84     st->st_blocks = partition->length_bytes / partition->disk_bytes_per_sector;
85     st->st_ctime = 0;
86     st->st_dev = 0xBAFEBABE;
87     st->st_gid = 0;
88     st->st_ino = 0xBAFEBABE;
89     st->st_mtime = 0;
90     st->st_nlink = 1;
91     st->st_rdev = 0xBAFEBABE;
92     st->st_size = partition->length_bytes;
93     st->st_uid = 0;
94     return 0;
95 }
96
97 static int Follow(struct vnode* node, struct vnode** out, const char* name) {
98     struct partition_data* partition = node->data;
99
100     if (!strcmp(name, "fs")) {
101         if (partition->fs == NULL) {
102             return EINVAL;
103         }
104
105         *out = partition->fs->node;
106         return 0;
107     }
108
109     return EINVAL;
110 }
111
112 static const struct vnode_operations dev_ops = {
113     .is_seekable = IsSeekable,

```

```

114 .read      = Read,
115 .write     = Write,
116 .create    = Create,
117 .follow    = Follow,
118 .dirent_type = DirentType,
119 .stat      = Stat,
120 ;
121
122 struct open_file* CreatePartition(struct open_file* disk, uint64_t start, uint64_t length, int id, int sector_size, int media_type, bool boot) {
123     struct partition_data* data = AllocHeap(sizeof(struct partition_data));
124     data->disk = disk;
125     data->disk_bytes_per_sector = sector_size;
126     data->id = id;
127     data->length_bytes = length;
128     data->start_byte = start;
129     data->media_type = media_type;
130     data->boot = boot;
131     data->fs = NULL;
132
133     struct vnode* node = CreateVnode(dev_ops);
134     node->data = data;
135
136     struct open_file* partition = CreateOpenFile(node, 0, 0, true, true);
137     LogWriteSerial("created the partition...\n");
138     MountFilesystemForDisk(partition);
139     return partition;
140 }
141
142 struct open_file* CreateMbrPartitionIfExists(struct open_file* disk, uint8_t* mem, int index, int sector_size) {
143     int offset = 0x1BE + index * 16;
144
145     uint8_t active = mem[offset + 0];
146     if (active & 0x7F) {
147         return NULL;
148     }
149
150     int media_type = mem[offset + 4];
151
152     uint32_t start_sector = mem[offset + 11];
153     start_sector <<= 8;
154     start_sector |= mem[offset + 10];
155     start_sector <<= 8;
156     start_sector |= mem[offset + 9];
157     start_sector <<= 8;
158     start_sector |= mem[offset + 8];
159
160     uint32_t total_sectors = mem[offset + 15];
161     total_sectors <<= 8;
162     total_sectors |= mem[offset + 14];
163     total_sectors <<= 8;
164     total_sectors |= mem[offset + 13];
165     total_sectors <<= 8;
166     total_sectors |= mem[offset + 12];
167
168     if (start_sector == 0 && total_sectors == 0) {
169         return NULL;
170     }
171
172     return CreatePartition(disk, ((uint64_t) start_sector) * sector_size, ((uint64_t) total_sectors) * sector_size, index, sector_size, media_type, active & 0x80);
173 }
174
175 /*
176  * caller to free return value.
177  */
178 struct open_file** GetMbrPartitions(struct open_file* disk) {
179     struct stat st;
180     int res = VnodeOpStat(disk->node, &st);
181     if (res != 0) {
182         return NULL;
183     }
184
185     uint8_t* mem = (uint8_t*) MapVirt(0, 0, st.st_blksize, VM_READ | VM_FILE, disk, 0);
186     if (mem == NULL) {
187         return NULL;
188     }
189
190     if (mem[0x1FE] != 0x55) {
191         return NULL;
192     }
193     if (mem[0x1FF] != 0xAA) {
194         return NULL;
195     }
196
197     struct open_file** partitions = AllocHeap(sizeof(struct open_file) * 5);
198     inline_memset(partitions, 0, sizeof(struct open_file) * 5);
199
200     int partitions_found = 0;
201     for (int i = 0; i < 4; ++i) {
202         struct open_file* partition = CreateMbrPartitionIfExists(disk, mem, i, st.st_blksize);
203         if (partition != NULL) {
204             partitions[partitions_found++] = partition;
205         }
206     }
207
208     UnmapVirt((size_t) mem, st.st_blksize);
209
210     return partitions;
211 }
212
213 /*
214  * null terminated array of struct vnode*
215  * e.g. {vnode_ptr_1, vnode_ptr_2, vnode_ptr_3, NULL}
216  */
217 struct open_file** GetPartitionsForDisk(struct open_file* disk) {
218     struct open_file** partitions = GetMbrPartitions(disk);
219
220     if (partitions == NULL) {
221         // check for GPT
222     }
223
224     return partitions;
225 }

```

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 }
21
22 return 0;
23 }
24
25 static uint8_t DirentType(struct vnode*) {
26     return DT_CHR;
27 }
28
29 static int Stat(struct vnode*, struct stat* st) {
30     st->st_mode = S_IFCHR | S_IRWXU | S_IRWXG | S_IRWXO;
31     st->st_atime = 0;
32     st->st_blksize = 0;
33     st->st_blocks = 0;
34     st->st_ctime = 0;
35     st->st_dev = 0xABEBCAFE;
36     st->st_gid = 0;
37     st->st_ino = 0xCAFEBAFE;
38     st->st_mtime = 0;
39     st->st_nlink = 1;
40     st->st_rdev = 0xCAFEDEAD;
41     st->st_size = 0;
42     st->st_uid = 0;
43     return 0;
44 }
45
46 static const struct vnode_operations dev_ops = {
47     .read = Read,
48     .dirent_type = DirentType,
49     .stat = Stat,
50 };
51
52 void InitRandomDevice(void)
53 {
54     AddVfsMount(CreateVnode(dev_ops), "rand");
55 }

```

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 pty_master_internal_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 pty_subordinate_internal_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 pty_master_internal_data* internal = node->data;
42     while (tr->length_remaining > 0) {
43         char c = BlockingBufferGet(internal->display_buffer);
44         PerformTransfer(&c, tr, 1);
45     }
46 }
47
48 return 0;
49 }
50
51 // "THE KEYBOARD"
52 static int MasterWrite(struct vnode* node, struct transfer* tr) {
53     struct pty_master_internal_data* internal = node->data;
54     while (tr->length_remaining > 0) {
55         char c;
56         PerformTransfer(&c, tr, 1);
57         BlockingBufferAdd(internal->keybrd_buffer, c, true);
58     }
59 }
60
61 return 0;
62 }
63
64 static int MasterWait(struct vnode*, int, uint64_t) {
65     return ENOSYS;
66 }
67
68 static int SubordinateWait(struct vnode*, int, uint64_t) {
69     return ENOSYS;
70 }

```

```

71 static void FlushSubordinateLineBuffer(struct vnode* node) {
72     struct pty_subordinate_internal_data* internal = node->data;
73     struct pty_master_internal_data* master_internal = internal->master->data;
74
75     // could add a 'BlockingBufferAddMany' call?
76     for (int i = 0; i < internal->line_buffer_pos; ++i) {
77         BlockingBufferAdd(master_internal->flushed_buffer, internal->line_buffer[i], true);
78     }
79
80     internal->line_buffer_pos = 0;
81 }
82
83 static void RemoveFromSubordinateLineBuffer(struct vnode* node) {
84     struct pty_subordinate_internal_data* internal = node->data;
85
86     if (internal->line_buffer_pos == 0) {
87         return;
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 pty_subordinate_internal_data* internal = node->data;
95
96     if (internal->line_buffer_pos == LINE_BUFFER_SIZE) {
97         Panic(PANIC_NOT_IMPLEMENTED);
98         return;
99     }
100
101     internal->line_buffer[internal->line_buffer_pos] = c;
102     internal->line_buffer_char_width[internal->line_buffer_pos] = width;
103     internal->line_buffer_pos++;
104 }
105
106 static void LineProcessor(void* sub) {
107     SetThreadPriority(GetThread(), SCHEDULE_POLICY_FIXED, FIXED_PRIORITY_KERNEL_HIGH);
108
109     struct vnode* node = (struct vnode*) sub;
110     struct pty_subordinate_internal_data* internal = node->data;
111     struct pty_master_internal_data* master_internal = internal->master->data;
112
113     while (true) {
114         bool echo = internal->termios.c_lflag & ECHO;
115         bool canon = internal->termios.c_lflag & ICANON;
116
117         char c = BlockingBufferGet(master_internal->keybrd_buffer);
118
119         /*
120          * This must happen before we modify the line buffer (i.e. to add or backspace a character), as
121          * the backspace code here needs to check for a non-empty line (and so this must be done before we make
122          * the line empty).
123          */
124         if (echo) {
125             if (c == '\b' && canon) {
126                 if (internal->line_buffer_pos > 0) {
127                     BlockingBufferAdd(master_internal->display_buffer, '\b', true);
128                     BlockingBufferAdd(master_internal->display_buffer, ' ', true);
129                     BlockingBufferAdd(master_internal->display_buffer, '\b', true);
130                 }
131             } else {
132                 BlockingBufferAdd(master_internal->display_buffer, c, true);
133             }
134         }
135
136         if (c == '\b' && canon) {
137             RemoveFromSubordinateLineBuffer(node);
138         }
139         else {
140             AddToSubordinateLineBuffer(node, c, 1);
141         }
142
143         if (c == '\n' || c == 3 || !canon) {
144             FlushSubordinateLineBuffer(node);
145         }
146     }
147 }
148
149 // "THE STDIN LINE BUFFER"
150 static int SubordinateRead(struct vnode* node, struct transfer* tr) {
151     struct pty_subordinate_internal_data* internal = (struct pty_subordinate_internal_data*) node->data;
152     struct pty_master_internal_data* master_internal = (struct pty_master_internal_data*) internal->master->data;
153
154     if (tr->length_remaining == 0) {
155         return 0;
156     }
157
158     char c = BlockingBufferGet(master_internal->flushed_buffer);
159     PerformTransfer(&c, tr, 1);
160
161     int res = 0;
162     while (tr->length_remaining > 0 && !(res = BlockingBufferTryGet(master_internal->flushed_buffer, (uint8_t*) &c))) {
163         PerformTransfer(&c, tr, 1);
164     }
165
166     return 0;
167 }
168
169 // "WRITING TO STDOUT"
170 static int SubordinateWrite(struct vnode* node, struct transfer* tr) {
171     struct pty_subordinate_internal_data* internal = (struct pty_subordinate_internal_data*) node->data;
172     struct pty_master_internal_data* master_internal = (struct pty_master_internal_data*) internal->master->data;
173
174     while (tr->length_remaining > 0) {
175         char c;
176         int err = PerformTransfer(&c, tr, 1);
177         if (err) {
178             return err;
179         }
180
181         BlockingBufferAdd(master_internal->display_buffer, c, true);
182     }
183
184     return 0;
185 }
186
187 static uint8_t GenericDirentType(struct vnode*) {
188     return DT_CHR;
189 }
190
191 static int GenericStat(struct vnode*, struct stat* st) {
192     st->st_mode = S_IFCHR | S_IRWXU | S_IRWXG | S_IRWXO;
193     st->st_atime = 0;
194     st->st_blksize = 0;
195     st->st_blocks = 0;
196     st->st_ctime = 0;
197     st->st_dev = 0xBABECAFEB;
198     st->st_gid = 0;
199     st->st_ino = 0xCAFEFEB;
200     st->st_mtime = 0;

```

```

201     st->st_nlink = 1;
202     st->st_rdev = 0xCAFEDEAD;
203     st->st_size = 0;
204     st->st_uid = 0;
205     return 0;
206 }
207
208 static int SubordinateCheckTty(struct vnode*) {
209     return 0;
210 }
211
212 static const struct vnode_operations master_operations = {
213     .read      = MasterRead,
214     .write     = MasterWrite,
215     .dirent_type = GenericDirentType,
216     .stat      = GenericStat,
217     .wait      = MasterWait,
218 };
219
220 static const struct vnode_operations subordinate_operations = {
221     .check_tty = SubordinateCheckTty,
222     .read      = SubordinateRead,
223     .write     = SubordinateWrite,
224     .dirent_type = GenericDirentType,
225     .stat      = GenericStat,
226     .wait      = SubordinateWait,
227 };
228
229 void CreatePseudoTerminal(struct vnode** master, struct vnode** subordinate) {
230     struct vnode* m = CreateVnode(master_operations);
231     struct vnode* s = CreateVnode(subordinate_operations);
232
233     struct pty_master_internal_data* m_data = AllocHeap(sizeof(struct pty_master_internal_data));
234     struct pty_subordinate_internal_data* s_data = AllocHeap(sizeof(struct pty_subordinate_internal_data));
235
236     m_data->subordinate = s;
237     m_data->display_buffer = BlockingBufferCreate(INTERNAL_BUFFER_SIZE);
238     m_data->keyboard_buffer = BlockingBufferCreate(INTERNAL_BUFFER_SIZE);
239     m_data->flushed_buffer = BlockingBufferCreate(FLUSHED_BUFFER_SIZE);
240     m_data->line_processing_thread = CreateThread(LineProcessor, (void*) s, GetVas(), "line processor");
241
242     s_data->master = m;
243     s_data->termios.c_lflag = ICANON | ECHO;
244
245     m->data = m_data;
246     s->data = s_data;
247     *master = m;
248     *subordinate = s;
249 }

```

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 Read(struct vnode*, struct transfer*) {
13     return 0;
14 }
15
16 static int Write(struct vnode*, struct transfer*) {
17     // TODO: do we need to set io->length_remaining to zero?
18     return 0;
19 }
20
21 static uint8_t DirentType(struct vnode*) {
22     return DT_CHR;
23 }
24
25 static int Stat(struct vnode*, struct stat* st) {
26     st->st_mode = S_IFCHR | S_IRWXU | S_IRWXG | S_IRWXO;
27     st->st_atime = 0;
28     st->st_blksize = 0;
29     st->st_blocks = 0;
30     st->st_ctime = 0;
31     st->st_dev = 0xBABECADE;
32     st->st_gid = 0;
33     st->st_ino = 0xCAFEBADE;
34     st->st_mtime = 0;
35     st->st_nlink = 1;
36     st->st_rdev = 0xCAFEDEAD;
37     st->st_size = 0;
38     st->st_uid = 0;
39     return 0;
40 }
41
42 static const struct vnode_operations dev_ops = {
43     .read      = Read,
44     .write     = Write,
45     .dirent_type = DirentType,
46     .stat      = Stat,
47 };
48
49 void InitNullDevice(void)
50 {
51     AddVfsMount(CreateVnode(dev_ops), "null");
52 }

```

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 irq) {
11     assert(strlen(name) <= 15);
12
13     if (irq < IRQL_SCHEDULER) {
14         Panic(PANIC_SPINLOCK_WRONG_IRQ);
15     }
16
17     lock->lock = 0;
18     lock->owner = NULL;
19     lock->irq = irq;
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->irq == 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_irq = GetIrql();
56     RaiseIrql(&lock->irq);
57
58     if (lock->irq != GetIrql()) {
59         Panic(PANIC_SPINLOCK_WRONG_IRQ);
60     }
61
62     AcquireSpinlockDirect(lock);
63     lock->prev_irq = prior_irq;
64     return prior_irq;
65 }
66
67 void ReleaseSpinlockIrql(struct spinlock* lock) {
68     int old_irq = lock->prev_irq;
69     ReleaseSpinlockDirect(lock);
70     LowerIrql(old_irq);
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  * @maxirq 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  *
55  * ETIMEDOUT if the semaphore was not acquired, and the operation timed out
56  *
57  * EAGAIN if the semaphore was not acquired, and the timeout_ms value was 0
58  *
59  * @maxirq1 IRQL_PAGE_FAULT
60  */
61 int AcquireSemaphore(struct semaphore* sem, int timeout_ms) {
62     MAX_IRQL(IRQL_PAGE_FAULT);
63     assert(sem != NULL);
64     LockScheduler();
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     return thr->timed_out ? (timeout_ms == 0 ? EAGAIN : ETIMEDOUT) : 0;
108 }
109
110 /**
111  * Releases (i.e., does the signal, or V operation on) a semaphore. If there are threads waiting on this semaphore,
112  * it will cause the first one to wake up.
113  *
114  * @param sem The semaphore to release/signal
115  */
116 void ReleaseSemaphore(struct semaphore* sem) {
117     MAX_IRQL(IRQL_PAGE_FAULT);
118     LockScheduler();
119     assert(sem->current_count > 0);
120
121     if (sem->waiting_list.head == NULL) {
122         if (sem->current_count == 0) {
123             Panic(PANIC_NEGATIVE_SEMAPHORE);
124         }
125         sem->current_count--;
126     } else {
127         struct thread* top = ThreadListDeleteTop(&sem->waiting_list);
128
129         /*
130          * If it's in the THREAD_STATE_WAITING_FOR_SEMAPHORE_WITH_TIMEOUT state, it could mean one of two things:
131          * - 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
132          * - it's been taken off the sleep queue and onto the ready already, but it hasn't yet been run yet (and is therefore
133          *   still in this state)
134          */
135         if (top->state == THREAD_STATE_WAITING_FOR_SEMAPHORE_WITH_TIMEOUT) {
136             bool on_sleep_queue = TryDequeueForSleep(top);
137
138             if (on_sleep_queue) {
139                 /*
140                  * Change the state to prevent UnblockThread from seeing it's in the timeout state and calling CancelSemaphoreOfThread.
141                  * If CancelSemaphoreOfThread were called, then it would attempt to delete it from the queue - but it's already been
142                  * deleted by this point and so would crash.
143                  */
144                 top->state = THREAD_STATE_READY;
145                 UnblockThread(top);
146             }
147             /*
148              * Do not unblock the thread if it's not on the sleep queue, as not being on the sleep queue means it's
149              * already on the ready queue.
150              */
151         } else {
152             UnblockThread(top);
153         }
154     }
155     UnlockScheduler();
156 }
157
158 /**
159  * Deallocates a semaphore. Loses its shit if someone is still holding onto it, as it is probably a bug if you're trying
160  * to destroy a semaphore when there's a possibility that someone might even be thinking about trying to acquire it (which
161  * would then try to acquire a deleted memory region, which is very bad).
162  *
163  * @param sem The semaphore to destroy.
164  * @param flags One of SEM_DONT_CARE, SEM_REQUIRE_ZERO or SEM_REQUIRE_FULL.
165  * @maxirq1 IRQL_SCHEDULER
166  */
167 int DestroySemaphore(struct semaphore* sem, int flags) {
168     MAX_IRQL(IRQL_SCHEDULER);
169     LockScheduler();
170     if (flags == SEM_REQUIRE_ZERO && sem->current_count != 0) {
171         UnlockScheduler();
172         return EBUSY;
173     }
174     if (flags == SEM_REQUIRE_FULL && sem->current_count != sem->max_count) {
175         UnlockScheduler();
176     }
177 }

```

```
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: ./merged.pdf

[binary]

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     LogWriteSerial("mounting filesystem for disk...\n");
62
63     AcquireMutex(fs_table_lock, -1);
64
65     struct open_file* fs = NULL;
66
67     for (int i = 0; i < num_filesystems; ++i) {
68         fs = NULL;
69         int res = registered_filesystems[i].mount_creator(partition, &fs);
70         if (res == 0) {
71             break;
72         }
73     }
74
75     ReleaseMutex(fs_table_lock);
76
77     if (fs == NULL) {
78         return ENODEV;
79     }
80
81     int res = VnodeOpCreate(partition->node, &fs->node, "fs", 0, 0);
82     if (res != 0) {
83         return res;
84     }
85
86     AddVfsMount(fs->node, GenerateNewMountedDiskName());
87     return 0;
88 }

```

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 | (IU << 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 /*
179 * Now we need to move on to the next entry if there is one.
180 */
181 if (buffer[0] == 0xFF) {
182     /* No more entries. */
183     return ENOENT;
184 }
185
186 } else if (buffer[0] == 0xFE) {
187     /*
188     * There is another entry, so read its inode and keep the loop going
189     */
190     parent = buffer[1];
191     parent |= (ino_t) buffer[2] << 8;
192     parent |= (ino_t) buffer[3] << 16;
193
194     /*
195     * Add the directory bit to the inode number as it should be a directory.
196     */
197     parent = INODE_TO_DIR(parent);
198 }
199
200 /*
201 * Something went very wrong if the directory header is not present!
202 */
203 return EIO;
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 bool IsSeekable(struct vnode*) {
39     return true;
40 }
41
42 static int CheckTty(struct vnode*) {
43     return ENOTTY;
44 }
45
46 static int Read(struct vnode* node, struct transfer* io) {
47     struct vnode_data* data = node->data;
48     if (data->directory)
49         return demofs_read_directory_entry(&data->fs, data->inode, io);
50     else {
51         return demofs_read_file(&data->fs, data->inode, data->file_length, io);
52     }
53 }
54
55 static int Write(struct vnode*, struct transfer*) {
56     return EROFS;
57 }
58
59 static int Create(struct vnode*, struct vnode**, const char*, int, mode_t) {
60     return EROFS;
61 }
62
63 static uint8_t DirentType(struct vnode* node) {
64     struct vnode_data* data = node->data;
65     return data->directory ? DT_DIR : DT_REG;
66 }
67
68 static int Stat(struct vnode* node, struct stat* stat) {
69     struct vnode_data* data = node->data;
70
71     stat->st_atime = 0;
72     stat->st_blksize = 512;
73     stat->st_blocks = 0;
74     stat->st_ctime = 0;
75     stat->st_dev = 0xDEADDEAD;
76     stat->st_gid = 0;
77     stat->st_ino = data->inode;
78     stat->st_mode = (INODE_IS_DIR(data->inode) ? S_IFDIR : S_IFREG) | S_IRWXU | S_IRWXG | S_IRWXO;
79     stat->st_mtime = 0;
80     stat->st_nlink = 1;
81     stat->st_rdev = 0;
82     stat->st_size = data->file_length;
83     stat->st_uid = 0;
84
85     return 0;

```

```

86 }
87
88 static int Truncate(struct vnode*, off_t) {
89     return EROFS;
90 }
91
92 static int Close(struct vnode* node) {
93     FreeHeap(node->data);
94     return 0;
95 }
96
97 static struct vnode* CreateDemoFsVnode();
98
99 static int Follow(struct vnode* node, struct vnode** out, const char* name) {
100     struct vnode_data* data = node->data;
101     if (data->directory) {
102         ino_t child_inode;
103         uint32_t file_length;
104
105         int status = demofs_follow(&data->fs, data->inode, &child_inode, name, &file_length);
106         if (status != 0) {
107             return status;
108         }
109
110         /*
111          * TODO: return existing vnode if someone opens the same file twice...
112          */
113
114         struct vnode* child_node = CreateDemoFsVnode();
115         struct vnode_data* child_data = AllocHeap(sizeof(struct vnode_data));
116         child_data->inode = child_inode;
117         child_data->fs = data->fs;
118         child_data->file_length = file_length;
119         child_data->directory = INODE_IS_DIR(child_inode);
120
121         child_node->data = child_data;
122
123         *out = child_node;
124
125         return 0;
126     }
127     else {
128         return ENOTDIR;
129     }
130 }
131
132 static const struct vnode_operations dev_ops = {
133     .check_open = CheckOpen,
134     .ioctl = Ioctl,
135     .is_seekable = IsSeekable,
136     .check_tty = CheckTty,
137     .read = Read,
138     .write = Write,
139     .close = Close,
140     .truncate = Truncate,
141     .create = Create,
142     .follow = Follow,
143     .dirent_type = DirentType,
144     .stat = Stat,
145 };
146
147 static struct vnode* CreateDemoFsVnode() {
148     return CreateVnode(dev_ops);
149 }
150
151 static int CheckForDemofsSignature(struct open_file* raw_device) {
152     struct stat st;
153     VnodeOpStat(raw_device->node, &st);
154
155     uint8_t* buffer = AllocHeapEx(st.st_blksize, HEAP_ALLOW_PAGING);
156     struct transfer io = CreateKernelTransfer(buffer, st.st_blksize, 8 * st.st_blksize, TRANSFER_READ);
157     int res = ReadFile(raw_device, &io);
158     if (res != 0) {
159         FreeHeap(buffer);
160         return ENOTSUP;
161     }
162
163     /*
164      * Check for the DemoFS signature.
165      */
166     if (buffer[0] != 'D' || buffer[1] != 'E' || buffer[2] != 'M' || buffer[3] != 'O') {
167         FreeHeap(buffer);
168         return ENOTSUP;
169     }
170
171     return 0;
172 }
173
174 int DemofsMountCreator(struct open_file* raw_device, struct open_file** out) {
175     int sig_check = CheckForDemofsSignature(raw_device);
176     if (sig_check != 0) {
177         return sig_check;
178     }
179
180     struct vnode* node = CreateDemoFsVnode();
181     struct vnode_data* data = AllocHeap(sizeof(struct vnode_data));
182
183     data->fs.disk = raw_device;
184     data->fs.root_inode = 9 | (1 << 31);
185     data->inode = 9 | (1 << 31); /* root directory inode */
186     data->file_length = 0; /* root directory has no length */
187     data->directory = true;
188
189     node->data = data;
190
191     *out = CreateOpenFile(node, 0, 0, true, false);
192     return 0;
193 }

```

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 <proglload.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),
51 * and then we can check if they have been modified. If they are, we will throw a
52 * somewhat nicer error than a system reboot.
53 *
54 * Note that we can still overflow 'badly' if someone makes an allocation on the
55 * stack lwhich is larger than the remaining space on the stack and the canary size
56 * combined.
57 *
58 * If the canary page is only partially used for the canary, the remainder of the
59 * page is able to be used normally.
60 */
61 #ifdef NDEBUG
62 #define NUM_CANARY_PAGES 0
63 #else
64 #define NUM_CANARY_BYTES (1024 * 8)
65 #define NUM_CANARY_PAGES BytesToPages(NUM_CANARY_BYTES)
66 #define CANARY_VALUE 0xBADF00D
67
68 static void CheckCanary(size_t canary_base) {
69     uint32_t* canary_ptr = (uint32_t*) canary_base;
70     for (size_t i = 0; i < NUM_CANARY_BYTES / sizeof(uint32_t); ++i) {
71         if (*canary_ptr++ != CANARY_VALUE) {
72             Panic(PANIC_CANARY_DIED);
73         }
74     }
75 }
76
77 static void CreateCanary(size_t canary_base) {
78     uint32_t* canary_ptr = (uint32_t*) canary_base;
79     for (size_t i = 0; i < NUM_CANARY_BYTES / sizeof(uint32_t); ++i) {
80         *canary_ptr++ = CANARY_VALUE;
81     }
82 }
83
84 #endif
85
86 /*
87 * Allocates a new page-aligned stack for a kernel thread, and returns
88 * the address of either the top of the stack (if it grows downward),
89 * or the bottom (if it grows upward).
90 */
91 static void CreateKernelStacks(struct thread* thr, int kernel_stack_kb) {
92     int total_bytes = (BytesToPages(kernel_stack_kb * 1024) + NUM_CANARY_PAGES) * ARCH_PAGE_SIZE;
93     size_t stack_bottom = MapVirt(0, 0, total_bytes, VM_READ | VM_WRITE | VM_LOCK, NULL, 0);
94     size_t stack_top = stack_bottom + total_bytes;
95
96 #ifdef NDEBUG
97     thr->canary_position = stack_bottom;
98     CreateCanary(stack_bottom);
99 #endif
100     thr->kernel_stack_top = stack_top;
101     thr->kernel_stack_size = total_bytes;
102     thr->stack_pointer = ArchPrepareStack(thr->kernel_stack_top);
103
104 static size_t CreateUserStack(int size) {
105     /*
106      * All user stacks share the same area of virtual memory, but have different
107      * mappings to physical memory.
108      */
109     int total_bytes = BytesToPages(size) * ARCH_PAGE_SIZE;
110     size_t stack_base = ARCH_USER_STACK_LIMIT - total_bytes;
111     size_t actual_base = MapVirt(0, stack_base, total_bytes, VM_READ | VM_WRITE | VM_USER | VM_LOCAL, NULL, 0);
112     assert(stack_base == actual_base);
113     (void) actual_base; // The assert gets taken out on release mode, so make the compiler happy
114     return ARCH_USER_STACK_LIMIT;
115 }
116
117 /**
118 * Blocks the currently executing thread (no effect will happen until the IRQ goes below IRQ_SCHEDULER).
119 * The scheduler lock must be held.
120 */
121 void BlockThread(int reason) {
122     AssertSchedulerLockHeld();
123     assert(reason != THREAD_STATE_READY && reason != THREAD_STATE_RUNNING);
124     assert(GetCpu() != NULL && GetCpu()->current_thread != NULL);
125     assert(GetCpu()->current_thread->state == THREAD_STATE_RUNNING);
126     GetCpu()->current_thread->state = reason;
127     PostponeScheduleUntilStandardIrql();
128 }
129
130 void UnblockThread(struct thread* thr) {
131     AssertSchedulerLockHeld();
132     if (thr->state == THREAD_STATE_WAITING_FOR_SEMAPHORE_WITH_TIMEOUT) {
133         CancelSemaphoreOfThread(thr);
134     }
135 }

```

```

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("COULD NOT 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
253     UpdateTimesliceExpiry();
254
255     /*
256      * To get here, someone must have called thread_schedule(), and therefore
257      * the lock must have been held.
258      */
259     UnlockScheduler();
260
261     /* Anything else you might want to do should be done here... */
262
263     /* Go to the address the thread actually wanted. */
264     GetThread()->initial_address(GetThread()->argument);
265
266     /* The thread has returned, so just terminate it. */
267     TerminateThread(GetThread());
268
269     Panic(PANIC_IMPOSSIBLE_RETURN);
270 }
271
272 static int GetMinPriorityValueForPolicy(int policy) {
273     if (policy == SCHEDULE_POLICY_USER_HIGHER) return 50;
274     if (policy == SCHEDULE_POLICY_USER_NORMAL) return 100;

```

```

276     if (policy == SCHEDULE_POLICY_USER_LOWER) return 150;
277     return 0;
278
279
280 static int GetMaxPriorityValueForPolicy(int policy) {
281     if (policy == SCHEDULE_POLICY_FIXED) return 255;
282     return GetMinPriorityValueForPolicy(policy) + 100;
283 }
284
285 static void UpdatePriority(bool yielded) {
286     struct thread* thr = GetThread();
287     int policy = thr->schedule_policy;
288
289     if (policy != SCHEDULE_POLICY_FIXED) {
290         int new_val = thr->priority + (yielded ? -1 : 1);
291         if (new_val >= GetMinPriorityValueForPolicy(policy) && new_val <= GetMaxPriorityValueForPolicy(policy)) {
292             thr->priority = new_val;
293         }
294     }
295 }
296
297 void LockSchedulerX(void) {
298     AcquireSpinlockIrql(&scheduler_lock);
299 }
300
301 void UnlockSchedulerX(void) {
302     ReleaseSpinlockIrql(&scheduler_lock);
303 }
304
305 void AssertSchedulerLockHeld(void) {
306     assert(IsSpinlockHeld(&scheduler_lock));
307 }
308
309 __attribute__((returns_twice)) static void SwitchToNewTask(struct thread* old_thread, struct thread* new_thread) {
310     new_thread->state = THREAD_STATE_RUNNING;
311     ThreadListDeleteTop(&ready_list);
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 #define MAX_IRQL_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->scheduler_policy = policy;
461     }
462     if (priority != -1) {
463         thr->priority = priority;
464     }
465     return 0;
466 }
467
468
469
470
471 void AssignThreadToCpu(void) {
472     // NO-OP UNTIL SMP IMPLEMENTED
473 }
474
475 void UnassignThreadToCpu(void) {
476     // NO-OP UNTIL SMP IMPLEMENTED
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 <irq.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_IRQL(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_IRQL(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     return false;
101 }
102
103 void HandleSleepWakeup(void* sys_time_ptr) {
104     MAX_IRQL(IRQL_PAGE_FAULT);
105
106     if (GetThread() == NULL) {
107         return;
108     }
109
110     LockScheduler();
111     if (sleep_wakeups_posted > 0) {
112         --sleep_wakeups_posted;
113     }
114
115     uint64_t system_time = *((uint64_t*) sys_time_ptr);
116
117     /*
118     * Wake up any sleeping tasks that need it.
119     */
120     while (PriorityQueueGetUsedSize(sleep_queue) > 0) {
121         struct priority_queue_result res = PriorityQueuePeek(sleep_queue);
122
123         /*
124         * Check if it needs waking.
125         */
126         if (res.priority <= system_time) {
127             struct thread* thr = *((struct thread**) res.data);
128             thr->timed_out = true;
129             PriorityQueuePop(sleep_queue);
130             UnblockThread(thr);
131         } else {
132             /*
133             * If this one doesn't need waking, none of the others will either.
134             */
135             break;
136         }
137     }
138
139     /*
140     * Check for any tasks that are asleep but on the overflow list. This is slow, but will
141     * only happen if we have more than 32 sleeping tasks, so it should normally not take any time.
142     */
143     struct thread* iter = sleep_overflow_list.head;
144     while (iter) {
145         if (iter->sleep_expiry <= system_time) {
146             ThreadListDelete(&sleep_overflow_list, iter);
147             iter->timed_out = true;
148             UnblockThread(iter);
149             iter = sleep_overflow_list.head; // restart, as list changed
150         } else {
151             iter = iter->next[NEXT_INDEX_SLEEP];
152         }
153     }
154
155     UnlockScheduler();
156 }
157
158 void SleepUntil(uint64_t system_time_ns) {
159     MAX_IRQL(IRQL_PAGE_FAULT);
160
161     if (system_time_ns < GetSystemTimer()) {
162         return;
163     }
164
165     LockScheduler();
166     GetThread()->sleep_expiry = system_time_ns;
167     QueueForSleep(GetThread());
168     BlockThread(THREAD_STATE_SLEEPING);
169     UnlockScheduler();
170
171     void SleepNano(uint64_t delta_ns) {
172         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     AcquireSpinlockIrql(&pid_lock);
55     pid_t pid = next_pid++;
56     ReleaseSpinlockIrql(&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;
80     dummy.pid = pid;
81     struct process_table_node* actual = AvlTreeGet(process_table, (void*) &dummy);
82     AvlTreeDelete(process_table, (void*) actual);
83     FreeHeap(actual); // this was allocated on 'InsertIntoProcessTable'
84
85     ReleaseMutex(&process_table_mutex);
86 }
87
88 void LockProcess(struct process* prcss) {
89     AcquireMutex(&prcss->lock, -1);
90 }
91
92 void UnlockProcess(struct process* prcss) {
93     ReleaseMutex(&prcss->lock);
94 }
95
96 void InitProcess(void) {
97     InitSpinlock(&pid_lock, "pid", IRQL_SCHEDULER);
98     process_table_mutex = CreateMutex("prcss table");
99     process_table = AvlTreeCreate();
100     AvlTreeSetComparator(process_table, ProcessTableComparator);
101 }
102
103 struct process* CreateProcess(pid_t parent_pid) {
104     EXACT_IRQL(IRQL_STANDARD);
105
106     struct process* prcss = AllocHeap(sizeof(struct process));
107
108     prcss->lock = CreateMutex("prcss");
109     prcss->vas = CreateVas();
110     prcss->parent = parent_pid;
111     prcss->children = AvlTreeCreate();
112     prcss->threads = AvlTreeCreate();
113     prcss->killed_children_semaphore = CreateSemaphore("killed children", SEM_BIG_NUMBER, SEM_BIG_NUMBER);
114     prcss->retv = 0;
115     prcss->terminated = false;
116     prcss->pid = InsertIntoProcessTable(prcss);
117     prcss->filedes_table = CreateFileDescriptorTable();

```

```

118
119     if (parent_pid != 0) {
120         struct process* parent = GetProcessFromPid(parent_pid);
121         LockProcess(parent);
122         AvlTreeInsert(parent->children, (void*) prcss);
123         UnlockProcess(parent);
124     }
125
126     return prcss;
127 }
128
129 void AddThreadToProcess(struct process* prcss, struct thread* thr) {
130     LockProcess(prcss);
131     AvlTreeInsert(prcss->threads, (void*) thr);
132     thr->process = prcss;
133     UnlockProcess(prcss);
134 }
135
136 struct process* ForkProcess(void) {
137     MAX_IRQL(IRQL_PAGE_FAULT);
138
139     LockProcess(GetProcess());
140
141     struct process* new_process = CreateProcess(GetProcess()->pid);
142     DestroyVas(new_process->vas);
143
144     // TODO: there are probably more things to copy over in the future (e.g. list of open file descriptors, etc.)
145     // the open files, etc.
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
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  * Changes the parent of a parentless process. Used to ensure the initial process can always reap orphaned
216  * processes.
217  */
218 static void AdoptOrphan(struct process* adopter, struct process* ophan) {
219     LockProcess(adopter);
220
221     ophan->parent = adopter->pid;
222     AvlTreeInsert(adopter->children, (void*) ophan);
223     ReleaseSemaphore(adopter->killed_children_semaphore);
224
225     UnlockProcess(adopter);
226 }
227
228 /**
229  * Recursively converts all child processes in a process' thread tree into zombie processes.
230  *
231  * @param node The subtree to start from. NULL is acceptable, and is the recursion base case.
232  */
233 static void RecursivelyMakeChildrenOrphans(struct avl_node* node) {
234     if (node == NULL) {
235         return;
236     }
237
238     RecursivelyMakeChildrenOrphans(AvlTreeGetLeft(node));
239     RecursivelyMakeChildrenOrphans(AvlTreeGetRight(node));
240     AdoptOrphan(GetProcessFromPid(1), AvlTreeGetData(node));
241 }
242
243 /**
244  * Recursively terminates all threads in a process' thread tree.
245  */
246
247

```

```

248 * @param node The subtree to start from. NULL is acceptable, and is the recursion base case.
249 */
250 static void RecursivelyKillRemainingThreads(struct avl_node* node) {
251     if (node == NULL) {
252         return;
253     }
254     RecursivelyKillRemainingThreads(AvlTreeGetLeft(node));
255     RecursivelyKillRemainingThreads(AvlTreeGetRight(node));
256     struct thread* victim = AvlTreeGetData(node);
257     if (victim->state != THREAD_STATE_TERMINATED && !victim->needs_termination) {
258         TerminateThread(victim);
259     }
260 }
261
262 /**
263  * Does all of the required operations to kill a process. This is run in its own thread, without an owning
264  * process, so that a process doesn't try to delete itself (and therefore delete its stack).
265  *
266  * @param arg The process to kill (needs to be cast to struct process*)
267  */
268 static void KillProcessHelper(void* arg) {
269     struct process* prcss = arg;
270     assert(GetProcess() == NULL);
271     assert(GetVas() != prcss->vas); // we should be on GetKernelVas()
272     RecursivelyKillRemainingThreads(AvlTreeGetRootNode(prcss->threads));
273     RecursivelyMakeChildrenOrphans(AvlTreeGetRootNode(prcss->children));
274     AvlTreeDestroy(prcss->threads);
275     AvlTreeDestroy(prcss->children);
276     DestroyVas(prcss->vas);
277     prcss->terminated = true;
278     if (prcss->parent == 0) {
279         ReapProcess(prcss);
280     } else {
281         struct process* parent = GetProcessFromPid(prcss->parent);
282         ReleaseSemaphore(parent->killed_children_semaphore);
283     }
284     TerminateThread(GetThread());
285 }
286
287 /**
288  * Deletes the current process and all its threads. Child processes have their parent switched to pid 1.
289  * If the process being deleted has a parent, then it becomes a zombie process until the parent reaps it
290  * If the process being deleted has no parent, it will be reaped and deallocated immediately.
291  * This function does not return.
292  *
293  * @param retv The return value the process being deleted will give.
294  */
295 void KillProcess(int retv) {
296     MAX_IRQL(IRQL_STANDARD);
297     struct process* prcss = GetProcess();
298     prcss->retv = retv;
299     /**
300      * Must run it in a different thread and process (a NULL process is fine), as it is going to kill all
301      * threads in the process, and the process itself.
302      */
303     CreateThreadEx(KillProcessHelper, (void*) prcss, GetKernelVas(), "process killer", NULL,
304                   SCHEDULE_POLICY_FIXED, FIXED_PRIORITY_KERNEL_HIGH, 0);
305     TerminateThread(GetThread());
306 }
307
308 /**
309  * Returns a pointer to the process that the thread currently running on this CPU belongs to. If there is no
310  * running thread (i.e. multitasking hasn't started yet), or the thread does not belong to a process, NULL
311  * is returned.
312  *
313  * @return The process of the current thread, if it exists, or NULL otherwise.
314  */
315 struct process* GetProcess(void) {
316     MAX_IRQL(IRQL_HIGH);
317     struct thread* thr = GetThread();
318     return thr == NULL ? NULL : thr->process;
319 }
320
321 struct process* CreateProcessWithEntryPoint(pid_t parent, void(*entry_point)(void*), void* args) {
322     EXACT_IRQL(IRQL_STANDARD);
323     struct process* prcss = CreateProcess(parent);
324     struct thread* thr = CreateThread(entry_point, args, prcss->vas, "prcssinit");
325     AddThreadToProcess(prcss, thr);
326     return prcss;
327 }
328
329 /**
330  * Returns the file descriptor table of the given process. Returns NULL if
331  * 'prcss' is null.
332  */
333 struct filedes_table* GetFileDescriptorTable(struct process* prcss) {
334     if (prcss == NULL) {
335         return NULL;
336     }
337     return prcss->filedes_table;
338 }
339
340 /**
341  * Given a process id, returns the process object. Returns NULL for an invalid
342  * 'pid'.
343  */
344 struct process* GetProcessFromPid(pid_t pid) {
345     EXACT_IRQL(IRQL_STANDARD);
346     AcquireMutex(process_table_mutex, -1);
347     struct process_table_node dummy;
348     dummy.pid = pid;
349     struct process_table_node* node = AvlTreeGet(process_table, (void*) &dummy);
350     ReleaseMutex(process_table_mutex);
351     return node == NULL ? NULL : node->process;
352 }
353
354 /**
355  * Returns the process id of a given process. If 'prcss' is null, 0 is returned.
356  */
357 pid_t GetPid(struct process* prcss) {
358     MAX_IRQL(IRQL_HIGH);

```



```

378     return prcss == NULL ? 0 : prcss->pid;
379 }
380
381 pid_t WaitProcess(pid_t pid, int* status, int flags) {
382     EXACT_IrqL(IrqL_STANDARD);
383
384     struct process* prcss = GetProcess();
385
386     pid_t result = 0;
387     int failed_reaps = 0;
388     while (result == 0) {
389         int res = AcquireSemaphore(prcss->killed_children_semaphore, (flags & WNOHANG) ? 0 : -1);
390         if (res != 0) {
391             break;
392         }
393         LockProcess(prcss);
394         result = RecursivelyTryReap(prcss, AvlTreeGetRootNode(prcss->children), pid, status);
395         UnlockProcess(prcss);
396         if (result == 0 && pid != (pid_t) -1) {
397             failed_reaps++;
398         }
399     }
400
401     /*
402      * Ensure that the next time we call WaitProcess(), we can immediately retry the reaps that
403      * we increased the semaphore for, but didn't actually reap on.
404      */
405     while (failed_reaps-- > 0) {
406         ReleaseSemaphore(prcss->killed_children_semaphore);
407     }
408
409     return result;
410 }

```

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     int res = OpenFile("sys:/progload.exe", O_RDONLY, 0, file);
24     if (res != 0) {
25         PanicEx(PANIC_PROGRAM_LOADER, "program loader couldn't be loaded");
26     }
27
28     res = GetFileSize(file, &program_loader_size);
29     if (res != 0) {
30         PanicEx(PANIC_PROGRAM_LOADER, "program loader size couldn't be found");
31     }
32
33     program_loader_addr = MapVirt(0, 0, program_loader_size, VM_READ | VM_FILE, file, 0);
34 }
35
36 int CopyProgramLoaderIntoAddressSpace(void) {
37     size_t mem = MapVirt(0, ARCH_PROG_LOADER_BASE, program_loader_size, VM_READ | VM_EXEC | VM_WRITE | VM_USER | VM_LOCAL, NULL, 0);
38     if (mem != ARCH_PROG_LOADER_BASE) {
39         return ENOMEM;
40     }
41
42     memcpy((void*) ARCH_PROG_LOADER_BASE, (void*) program_loader_addr, program_loader_size);
43     return 0;
44 }

```

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() {
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     int res = OpenFile("swap:/", O_RDWR, 0, &swapfile);
35     if (res != 0) {
36         Panic(PANIC_NO_FILESYSTEM);
37     }
38
39     InitSpinlock(&swapfile_lock, "swapfile", IRQL_SCHEDULER);
40     SetupSwapfileBitmap();
41 }
42
43 struct open_file* GetSwapfile(void) {
44     return swapfile;
45 }
46
47 static bool GetBitmapEntry(size_t index) {
48     return swapfile_bitmap[index / 8] & (1 << (index % 8));
49 }
50
51 static void SetBitmapEntry(size_t index, bool value) {
52     if (value) {
53         swapfile_bitmap[index / 8] |= 1 << (index % 8);
54     } else {
55         swapfile_bitmap[index / 8] &= ~(1 << (index % 8));
56     }
57 }
58
59 uint64_t AllocateSwapfileIndex(void) {
60     AcquireSpinlockIrql(&swapfile_lock);
61     ++number_on_swapfile;
62
63     for (size_t i = 0; i < num_swapfile_bitmap_entries; ++i) {
64         if (!GetBitmapEntry(i)) {
65             SetBitmapEntry(i, true);
66             ReleaseSpinlockIrql(&swapfile_lock);
67             return i;
68         }
69     }
70
71     Panic(PANIC_OUT_OF_SWAPFILE);
72 }
73
74 void DeallocateSwapfileIndex(uint64_t index) {
75     AcquireSpinlockIrql(&swapfile_lock);
76     --number_on_swapfile;
77     SetBitmapEntry(index, false);
78     ReleaseSpinlockIrql(&swapfile_lock);
79 }
80
81 int GetNumberOfPagesOnSwapfile(void) {
82     AcquireSpinlockIrql(&swapfile_lock);
83     int val = number_on_swapfile;
84     ReleaseSpinlockIrql(&swapfile_lock);
85     return val;
86 }

```

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 /*
17  * TODO: once semaphores are allowed / virtual memory is initialised, we need
18  * to swap from using spinlocks to using a mutex... the reason being that we
19  * can't actually manipulate the pageable heap while a spinlock is held! (as
20  * the block might be paged out!!)
21  */
22
23
24 #define BOOTSTRAP_AREA_SIZE (1024 * 16)
25 #define MAX_EMERGENCY_BLOCKS 16
26
27 /*
28  * For requests larger than this, we'll issue a warning to say that MapVirt is a much better choice.
29  */
30 #define WARNING_LARGE_REQUEST_SIZE (1024 * 3 + 512)
31
32 struct emergency_block {
33     uint8_t* address;
34     size_t size;
35     bool valid;
36 };
37

```

```

38 /**
39  * Used as a system block that can be used even before the virtual memory manager is setup.
40  */
41 static uint8_t bootstrap_memory_area[BOOTSTRAP_AREA_SIZE] __attribute__((aligned(ARCH_PAGE_SIZE)));
42
43 /**
44  * Used to give us memory when we are not allowed to fault (i.e. can't allocate virtual memory).
45  */
46 static struct emergency_block emergency_blocks[MAX_EMERGENCY_BLOCKS] = {
47     { .address = bootstrap_memory_area, .size = BOOTSTRAP_AREA_SIZE, .valid = true };
48 };
49
50 static struct semaphore heap_lock;
51 static struct spinlock heap_spinlock;
52
53 bool IsHeapReinitialised(void) {
54     return heap_lock != NULL;
55 }
56
57 static struct spinlock heap_locker_lock;
58 static struct thread* lock_entry_threads[2];
59
60 static bool LockHeap(bool paging) {
61     bool acquired = false;
62     struct thread* thr = GetThread();
63
64     AcquireSpinlockIrql(&heap_locker_lock);
65     if (lock_entry_threads[paging ? 1 : 0] != thr || lock_entry_threads[paging ? 1 : 0] == NULL) {
66         ReleaseSpinlockIrql(&heap_locker_lock);
67
68         // TODO: we have atomicity issues between the release of the previous lock
69         // and the acquisition of the next one... may need a retry loop
70
71         /*
72          * e.g.
73          * retry:
74          * Acquire(locker_lock);
75          * if (big long condition) {
76          *     if (paging) {
77          *         TryAcquire...
78          *     } else {
79          *         TryAcquire...
80          *     }
81          *
82          * if (acquired) {
83          *     if (!(the same big long condition)) {
84          *         release
85          *         goto retry;
86          *     }
87          *     goto retry;
88          * }
89          */
90
91         if (paging) {
92             AcquireMutex(&heap_lock, -1);
93         } else {
94             AcquireSpinlockIrql(&heap_spinlock);
95         }
96
97         // e.g. check that
98         // (lock_entry_threads[paging ? 1 : 0] != thr || lock_entry_threads[paging ? 1 : 0] == NULL)
99         // still holds,
100
101         acquired = true;
102         lock_entry_threads[paging ? 1 : 0] = thr;
103     }
104
105     return acquired;
106 }
107
108 static void UnlockHeap(bool paging) {
109     lock_entry_threads[paging ? 1 : 0] = NULL;
110
111     if (paging) {
112         ReleaseMutex(&heap_lock);
113     } else {
114         ReleaseSpinlockIrql(&heap_spinlock);
115     }
116 }
117
118 static void* AllocateFromEmergencyBlocks(size_t size) {
119     int smallest_block = -1;
120
121     for (int i = 0; i < MAX_EMERGENCY_BLOCKS; ++i) {
122         if (emergency_blocks[i].valid && emergency_blocks[i].size >= size) {
123             if (smallest_block == -1 || emergency_blocks[i].size < emergency_blocks[smallest_block].size) {
124                 smallest_block = i;
125             }
126         }
127     }
128
129     if (smallest_block == -1) {
130         Panic(PANIC_CANNOT_MALLOC_WITHOUT_FAULTING);
131     }
132
133     void* address = emergency_blocks[smallest_block].address;
134
135     emergency_blocks[smallest_block].address += size;
136     emergency_blocks[smallest_block].size -= size;
137
138     if (emergency_blocks[smallest_block].size < ARCH_PAGE_SIZE) {
139         emergency_blocks[smallest_block].valid = false;
140     }
141
142     return address;
143 }
144
145 /** This function needs to be called with the heap lock held.
146  */
147 static void AddBlockToBackupHeap(size_t size) {
148     UnlockHeap(false);
149     void* address = (void*) MapVirt(0, 0, size, VM_READ | VM_WRITE | VM_LOCK, NULL, 0);
150     LockHeap(false);
151
152     int index_of_smallest_block = 0;
153     for (int i = 0; i < MAX_EMERGENCY_BLOCKS; ++i) {
154         if (emergency_blocks[i].valid) {
155             if (emergency_blocks[i].size < emergency_blocks[index_of_smallest_block].size) {
156                 index_of_smallest_block = i;
157             }
158         } else {
159             emergency_blocks[i].valid = true;
160             emergency_blocks[i].size = size;
161             emergency_blocks[i].address = address;
162             return;
163         }
164     }
165 }

```

```

168     }
169 }
170
171 LogDeveloperWarning("losing 0x%X bytes due to strangeness with backup heap.\n", emergency_blocks[index_of_smallest_block].size);
172 LogDeveloperWarning("TODO: could probably just add this as a regular block to the regular heap.\n");
173
174 emergency_blocks[index_of_smallest_block].size = size;
175 emergency_blocks[index_of_smallest_block].address = address;
176
177
178 static void RestoreEmergencyPages(void* context) {
179     (void) context;
180
181     EXACT_IRQL(IRQL_STANDARD);
182
183     /*
184      * TODO: maybe make this greedier (i.e. grab larger blocks), but also have a way for the PMM to ask for larger
185      * blocks back if needed. would still need to retain enough stashed away for PMM/VMM to use the heap, and would
186      * need to unlock the pages, and ensure that allocations from emergency blocks are done via smallest-fit (so large
187      * blocks aren't wasted).
188      */
189
190     size_t total_size = 0;
191     size_t largest_block = 0;
192
193     LockHeap(false);
194     for (int i = 0; i < MAX_EMERGENCY_BLOCKS; ++i) {
195         if (emergency_blocks[i].valid) {
196             size_t size = emergency_blocks[i].size;
197             total_size += size;
198             if (size > largest_block) {
199                 largest_block = size;
200             }
201         }
202     }
203
204     while (largest_block < BOOTSTRAP_AREA_SIZE / 2 || total_size < BOOTSTRAP_AREA_SIZE) {
205         AddBlockToBackupHeap(BOOTSTRAP_AREA_SIZE);
206         total_size += BOOTSTRAP_AREA_SIZE;
207         largest_block = BOOTSTRAP_AREA_SIZE;
208     }
209
210     UnlockHeap(false);
211 }
212
213 /**
214  * Represents a section of memory that is either allocated or free. The memory address
215  * it represents is itself, excluding the metadata at the start or end.
216  *
217  * See the report for further details.
218  */
219 struct block {
220     /*
221      * The entire size of the block. The low 2 bits do not form part of the size,
222      * the low bit is set for allocated blocks, and the second lowest bit indicates
223      * if it is on the swappable heap or not.
224      */
225     size_t size;
226
227     /*
228      * Only here on free blocks. Allocated blocks use this as the start of allocated
229      * memory.
230      */
231     struct block* next;
232     struct block* prev;
233
234     /*
235      * At position size - sizeof(size_t), there is the trailing size tag.
236      * There are no flags in the low bit of this value, unlike the heading size tag.
237      */
238 };
239
240 #ifndef NDEBUG
241 static int outstanding_allocations = 0;
242
243 int DbgGetOutstandingHeapAllocations(void) {
244     return outstanding_allocations;
245 }
246 #endif
247
248 /**
249  * Must be a power of 2.
250  */
251 #define ALIGNMENT 8
252
253 /**
254  * The amount of metadata at the start and end of allocated blocks. The next and free
255  * pointers in free blocks do not count.
256  */
257 #define METADATA_LEADING_AMOUNT (sizeof(size_t))
258 #define METADATA_TRIALING_AMOUNT (sizeof(size_t))
259 #define METADATA_TOTAL_AMOUNT (METADATA_LEADING_AMOUNT + METADATA_TRIALING_AMOUNT)
260
261 /**
262  * (THIS COMMENT IS ABOUT x86-64 - halve the byte values for x86)
263  * Having blocks of size 8 is wasteful, as they need to be at least 32 bytes long total to fit
264  * the metadata when free, but only need 16 bytes metadata when allocated. Therefore, we have
265  * 8 spare bytes that are wasted.
266  */
267 #define MINIMUM_REQUEST_SIZE_INTERNAL (2 * sizeof(size_t))
268
269 /**
270  * The last size should be larger than the max allocation size, as otherwise we could allocate larger
271  * than we have in the final list.
272  */
273 #define TOTAL_NUM_FREE_LISTS 35
274
275 /**
276  * An array which holds the minimum allocation sizes that each free list can hold.
277  */
278 static size_t free_list_block_sizes[TOTAL_NUM_FREE_LISTS] = {
279     8,          16,          24,          32,
280     40,         48,          56,          64,
281     72,         80,          88,          96,
282     104,        112,         120,         128,
283     160,        192,         224,         256,
284     320,        384,         448,         512,
285     768,        1024,        1536,        2048,    // 28 [27]
286     1024 * 4,   1024 * 8,   1024 * 16,  1024 * 32,    // 32 [31]
287     1024 * 64,  1024 * 128, 1024 * 256,
288 };
289
290 /**
291  * Used to work out which free list a block should be in, when we are *reading* a block.
292  * This rounds the size *up*, meaning it cannot be used to insert a block into a list.
293  * NOT USED TO INSERT BLOCKS!!
294  */
295 static int GetSmallestListIndexThatFits(size_t size_without_metadata) {
296     int i = 0;
297     while (true) {

```

```

298     if (size_without_metadata <= free_list_block_sizes[i]) {
299         return i;
300     }
301     ++i;
302 }
303
304 /**
305  * Calculates which free list a block should be inserted in. This one rounds *down*, and
306  * so it should not normally be used to look up where a block should be.
307  */
308 static int GetInsertionIndex(size_t size_without_metadata) {
309     /*
310      * This will only fail if we somehow get a block that is smaller than the minimum
311      * possible size (i.e., something has gone very wrong.)
312      */
313     assert(size_without_metadata >= free_list_block_sizes[0]);
314
315     /*
316      * We can't round down to the next one when it's the smallest possible size,
317      * so handle this case specially.
318      */
319     if (size_without_metadata == free_list_block_sizes[0]) {
320         return 0;
321     }
322
323     /*
324      * Look until we go past the block size we need, and then return the previous
325      * value. This gives us the final block size that doesn't exceed the input value.
326      */
327     for (int i = 0; i < TOTAL_NUM_FREE_LISTS - 1; ++i) {
328         if (size_without_metadata <= free_list_block_sizes[i]) {
329             return i - 1;
330         }
331     }
332
333     Panic(PANIC_HEAP_REQUEST_TOO_LARGE);
334 }
335
336 /**
337  * Rounds up a user-supplied allocation size to the alignment. If the value is smaller than
338  * the minimum request size that is internally supported, it will round it up to that size.
339  */
340 static size_t RoundUpSize(size_t size) {
341     /* Should have already been checked against. */
342     assert(size != 0);
343
344     if (size < MINIMUM_REQUEST_SIZE_INTERNAL) {
345         size = MINIMUM_REQUEST_SIZE_INTERNAL;
346     }
347
348     return (size + ALIGNMENT - 1) & ~(ALIGNMENT - 1);
349 }
350
351 /**
352  * Marks a block as being free (unallocated).
353  */
354 static void MarkFree(struct block* block) {
355     block->size &= ~1;
356 }
357
358 /**
359  * Marks a block as being allocated.
360  */
361 static void MarkAllocated(struct block* block) {
362     block->size |= 1;
363 }
364
365 /**
366  * Returns true if the block is allocated, or false if it is free.
367  */
368 static bool IsAllocated(struct block* block) {
369     return block->size & 1;
370 }
371
372 static void MarkSwappability(struct block* block, int can_swap) {
373     if (can_swap) {
374         block->size |= 2;
375     } else {
376         block->size &= ~2;
377     }
378 }
379
380 static bool IsOnSwappableHeap(struct block* block) {
381     return block->size & 2;
382 }
383
384 /**
385  * Global arrays always initialise to zero (and therefore, to NULL).
386  * Entries in free lists must have a user allocated size GREATER OR EQUAL TO the size in free_list_block_sizes.
387  */
388 static struct block* _head_block[TOTAL_NUM_FREE_LISTS];
389 static struct block* _head_block_swappable[TOTAL_NUM_FREE_LISTS];
390
391 static struct block** GetHeap(bool swappable) {
392     return swappable ? _head_block_swappable : _head_block;
393 }
394
395 static struct block** GetHeapForBlock(struct block* block) {
396     return GetHeap(IsOnSwappableHeap(block));
397 }
398
399 /**
400  * Given a block, returns its total size, including metadata. This takes into account
401  * the flags on the size field and removes them from the return value.
402  */
403 static size_t GetBlockSize(struct block* block) {
404     size_t size = block->size & ~3;
405
406     /*
407      * Ensure the other size tag matches. If it doesn't, there has been memory corruption.
408      */
409     assert((((size_t*) block) + (size / sizeof(size_t)) - 1) == size);
410
411     return size;
412 }
413
414 void DbgPrintListStats(void) {
415     LogWriteSerial("\nDbgPrintListStats:\n");
416     for (int i = 0; i < TOTAL_NUM_FREE_LISTS; ++i) {
417         struct block* unswap = GetHeap[false][i];
418         struct block* swap = GetHeap[true][i];
419
420         size_t unswap_blocks = 0;
421         size_t unswap_size = 0;
422         size_t swap_blocks = 0;
423         size_t swap_size = 0;
424
425         int timeout = 0;
426     }
427 }

```

```

428     while (unswap) {
429         ++unswap_blocks;
430         unswap_size += GetBlockSize(unswap);
431         unswap = unswap->next;
432         if (timeout <= 100000) {
433             Panic(PANIC_DOUBLE_FREE_DETECTED);
434         }
435     }
436
437     timeout = 0;
438     while (swap) {
439         ++swap_blocks;
440         swap_size += GetBlockSize(swap);
441         swap = swap->next;
442         if (timeout >= 100000) {
443             Panic(PANIC_DOUBLE_FREE_DETECTED);
444         }
445     }
446
447     if ((unswap_blocks | swap_blocks) != 0) {
448         LogWriteSerial("    Bucket %d [0x%X]: unswappable %d / 0x%X. swappable %d / 0x%X\n", i, free_list_block_sizes[i], unswap_blocks, unswap_size, swap_bl
449     }
450 }
451 LogWriteSerial("\n\n");
452 }
453
454 /**
455  * Sets the *total* size of a given block. This does not do any checking, so the caller must be
456  * careful to ensure that calling this doesn't leak memory (by setting it too small), or cause
457  * double-allocation of the same area (by setting it too large). Carries the swappability and allocation
458  * tags with it from the front to the back tags.
459  */
460 static void SetSizeTags(struct block* block, size_t size) {
461     block->size = (block->size & 3) | size;
462     *((size_t*) block) + (size / sizeof(size_t) - 1) = size;
463 }
464
465 static void* GetSystemMemory(size_t size, int flags) {
466     EXACT_IQRL(IQRL_SCHEDULER);
467
468     if (flags & HEAP_NO_FAULT) {
469         if (flags & HEAP_ALLOW_PAGING) {
470             Panic(PANIC_CONFLICTING_ALLOCATION_REQUIREMENTS);
471         }
472         if (!IsVirtInitialised()) {
473             DeferUntilIrql(IQRL_STANDARD, RestoreEmergencyPages, NULL);
474         }
475         return AllocateFromEmergencyBlocks(size);
476     }
477
478     if (flags & HEAP_ALLOW_PAGING) {
479         LogWriteSerial("about to allocate %d KBs of pageable memory\n", (size + 4095) / 4096 * 4);
480     }
481     return (void*) MapVirt(0, 0, size, VM_READ | VM_WRITE | (flags & HEAP_ALLOW_PAGING ? 0 : VM_LOCK), NULL, 0);
482 }
483
484 /**
485  * Allocates a new block from the system that is able to hold the amount of data
486  * specified. Also allocated enough memory for fenceposts on either side of the data,
487  * and sets up these fenceposts correctly.
488  */
489 static struct block* RequestBlock(size_t total_size, int flags) {
490     EXACT_IQRL(IQRL_SCHEDULER);
491
492     /*
493      * We need to add the extra bytes for fenceposts to be added. We must do this before we
494      * round up to the nearest areana size (if we did it after, it wouldn't be aligned anymore).
495      */
496     total_size += MINIMUM_REQUEST_SIZE_INTERNAL * 2;
497     total_size = (total_size + ARCH_PAGE_SIZE - 1) & ~(ARCH_PAGE_SIZE - 1);
498
499     if (!IsVirtInitialised()) {
500         flags |= HEAP_NO_FAULT;
501     }
502
503     /*
504      * Get memory from the system.
505      */
506     struct block* block = (struct block*) GetSystemMemory(total_size, flags);
507     if (block == NULL) {
508         Panic(PANIC_OUT_OF_HEAP);
509     }
510
511     /*
512      * Set the metadata for both the fenceposts and the main data block.
513      * Keep in mind that total_size now includes the fencepost metadata (see top of function), so this
514      * sometimes needs to be subtracted off.
515      */
516     struct block* left_fence = block;
517     struct block* actual_block = ((struct block*) (((size_t*) block) + MINIMUM_REQUEST_SIZE_INTERNAL / sizeof(size_t)));
518     struct block* right_fence = ((struct block*) (((size_t*) block) + (total_size - MINIMUM_REQUEST_SIZE_INTERNAL) / sizeof(size_t)));
519
520     SetSizeTags(left_fence, MINIMUM_REQUEST_SIZE_INTERNAL);
521     SetSizeTags(actual_block, total_size - 2 * MINIMUM_REQUEST_SIZE_INTERNAL);
522     SetSizeTags(right_fence, MINIMUM_REQUEST_SIZE_INTERNAL);
523
524     actual_block->prev = NULL;
525     actual_block->next = NULL;
526
527     MarkAllocated(left_fence);
528     MarkAllocated(right_fence);
529     MarkFree(actual_block);
530
531     MarkSwappability(left_fence, flags & HEAP_ALLOW_PAGING);
532     MarkSwappability(right_fence, flags & HEAP_ALLOW_PAGING);
533     MarkSwappability(actual_block, flags & HEAP_ALLOW_PAGING);
534
535     return actual_block;
536 }
537
538 /**
539  * Removes a block from a free list. It needs to take in the exact free list's index (as opposed to calculating
540  * it itself), as this may be used halfway through allocations or deallocations where the block isn't yet in
541  * its correct block.
542  */
543 static void RemoveBlock(int free_list_index, struct block* block) {
544     EXACT_IQRL(IQRL_SCHEDULER);
545
546     struct block** head_list = GetHeapForBlock(block);
547
548     /*
549      * Perform a standard linked list deletion.
550      */
551     if (block->prev == NULL && block->next == NULL) {
552         assert(head_list[free_list_index] == block);
553         head_list[free_list_index] = NULL;
554     } else if (block->prev == NULL) {
555         head_list[free_list_index] = block->next;
556         block->next->prev = NULL;
557     }

```

```

558
559     } else if (block->next == NULL) {
560         block->prev->next = NULL;
561     }
562     } else {
563         block->prev->next = block->next;
564         block->next->prev = block->prev;
565     }
566 }
567
568 /**
569  * Adds a block to its appropriate free list. It also coalesces the block with surrounding free blocks
570  * if possible.
571  */
572 static struct block* AddBlock(struct block* block) {
573     EXACT_IRQL(IRQL_SCHEDULER);
574
575     /*
576      * Although this function is technically recursive (because it needs to shuffle blocks into different
577      * lists by calling itself again), but there are a constant number of free lists, so it is still
578      * coalescing in constant time.
579      */
580     size_t size = GetBlockSize(block);
581     struct block** head_list = GetHeapForBlock(block);
582
583     int free_list_index = GetInsertionIndex(size - METADATA_TOTAL_AMOUNT);
584
585     size_t prev_block_size = (((size_t*) block) - 1);
586     struct block* prev_block = (struct block*) (((size_t*) block) - prev_block_size / sizeof(size_t));
587     struct block* next_block = (struct block*) (((size_t*) block) + size / sizeof(size_t));
588
589     if (IsAllocated(prev_block) && IsAllocated(next_block)) {
590         /*
591          * Cannot coalesce here, so just add to the free list.
592          */
593         block->prev = NULL;
594         block->next = head_list[free_list_index];
595         if (block->next != NULL)
596             block->next->prev = block;
597
598         head_list[free_list_index] = block;
599         MarkFree(block);
600         return block;
601     }
602     } else if (IsAllocated(prev_block) && !IsAllocated(next_block)) {
603         /*
604          * Need to coalesce with the one on the right.
605          */
606
607         /*
608          * Swappable and non-swappable blocks should be on entirely separate heaps, and you can't look into
609          * the other because the fences should prevent anyone looking between them.
610          */
611         bool swappable = IsOnSwappableHeap(block);
612         assert(swappable == IsOnSwappableHeap(next_block));
613
614         RemoveBlock(GetInsertionIndex(GetBlockSize(next_block) - METADATA_TOTAL_AMOUNT), next_block);
615         SetSizeTags(block, size + GetBlockSize(next_block));
616         block->prev = NULL;
617         block->next = NULL;
618         MarkFree(block);
619         MarkSwappability(block, swappable);
620
621         return AddBlock(block);
622     }
623     } else if (!IsAllocated(prev_block) && IsAllocated(next_block)) {
624         /*
625          * Need to coalesce with the one on the left.
626          */
627
628         /*
629          * Swappable and non-swappable blocks should be on entirely separate heaps, and you can't look into
630          * the other because the fences should prevent anyone looking between them.
631          */
632         bool swappable = IsOnSwappableHeap(block);
633         assert(swappable == IsOnSwappableHeap(prev_block));
634
635         RemoveBlock(GetInsertionIndex(GetBlockSize(prev_block) - METADATA_TOTAL_AMOUNT), prev_block);
636         SetSizeTags(prev_block, size + GetBlockSize(prev_block));
637         prev_block->prev = NULL;
638         prev_block->next = NULL;
639         MarkFree(prev_block);
640         MarkSwappability(prev_block, swappable);
641
642         return AddBlock(prev_block);
643     }
644     } else {
645         /*
646          * Coalesce with blocks on both sides.
647          */
648
649         /*
650          * Swappable and non-swappable blocks should be on entirely separate heaps, and you can't look into
651          * the other because the fences should prevent anyone looking between them.
652          */
653         bool swappable = IsOnSwappableHeap(block);
654         assert(swappable == IsOnSwappableHeap(prev_block));
655         assert(swappable == IsOnSwappableHeap(next_block));
656
657         RemoveBlock(GetInsertionIndex(GetBlockSize(prev_block) - METADATA_TOTAL_AMOUNT), prev_block);
658         RemoveBlock(GetInsertionIndex(GetBlockSize(next_block) - METADATA_TOTAL_AMOUNT), next_block);
659
660         SetSizeTags(prev_block, size + GetBlockSize(prev_block) + GetBlockSize(next_block));
661         prev_block->prev = NULL;
662         prev_block->next = NULL;
663         MarkFree(prev_block);
664         MarkSwappability(prev_block, swappable);
665
666         return AddBlock(prev_block);
667     }
668 }
669
670 /**
671  * Allocates a block. The block to be allocated will be the first block in the given free
672  * list, and that free list must be non-empty, and be able to fit the requested size.
673  */
674 static struct block* AllocateBlock(struct block* block, int free_list_index, size_t user_requested_size) {
675     EXACT_IRQL(IRQL_SCHEDULER);
676     assert(block != NULL);
677
678     size_t total_size = user_requested_size + METADATA_TOTAL_AMOUNT;
679     size_t block_size = GetBlockSize(block);
680
681     assert(block_size >= total_size);
682
683     if (block_size - total_size < MINIMUM_REQUEST_SIZE_INTERNAL + METADATA_TOTAL_AMOUNT) {
684         /*
685          * We can just remove from the list altogether if the sizes match up exactly,
686          * or if there would be so little left over that we can't form a new block.
687          */

```



```

688     RemoveBlock(free_list_index, block);
689     /*
690     * Prevent memory leak (from having a hole in memory), but do it after removing
691     * the block, as this may change the list it needs to be in, and RemoveBlock
692     * will not like that.
693     */
694     SetSizeTags(block, block_size);
695     MarkAllocated(block);
696     return block;
697 } else {
698     /*
699     * We must split the block into two. If no list change is needed, we can leave the 'leftover' parts in the list
700     * as is (just fixing up the size tags), and then return the new block.
701     */
702     RemoveBlock(free_list_index, block);
703
704     size_t leftover = block_size - total_size;
705     SetSizeTags(block, leftover);
706
707     struct block* allocated_block = (struct block*) (((size_t*) block) + (leftover / sizeof(size_t)));
708     SetSizeTags(allocated_block, total_size);
709
710     /*
711     * We are giving it new tags, so must set this correctly.
712     */
713     MarkSwappability(allocated_block, IsOnSwappableHeap(block));
714
715     /*
716     * Must be done before we try to move around the leftovers (or else it will actually
717     * coalesce back into one block).
718     */
719     MarkAllocated(allocated_block);
720
721     /*
722     * We need to remove the leftover block from this list, and add it to the correct list.
723     */
724     MarkFree(block);
725     AddBlock(block);
726     return allocated_block;
727 }
728 }
729 }
730 }
731
732 /**
733  * Allocates a block that can fit the user requested size. It will request new memory from the
734  * system if required. If it returns NULL, then there is not enough memory of the system to
735  * satisfy the request.
736  */
737 static struct block* FindBlock(size_t user_requested_size, int flags) {
738     EXACT_IRQL(IRQL_SCHEDULER);
739
740     struct block** head_list = GetHeap(flags & HEAP_ALLOW_PAGING);
741     /*
742     * Check the free lists in order, starting from the smallest one that will fit the block.
743     * If we find one with a free block, then we allocate the head of that list.
744     */
745     int min_index = GetSmallestListIndexThatFits(user_requested_size);
746
747     for (int i = min_index; i < TOTAL_NUM_FREE_LISTS; ++i) {
748         if (head_list[i] != NULL) {
749             return AllocateBlock(head_list[i], i, user_requested_size);
750         }
751     }
752
753     /*
754     * If we can't find a block that will fit, then we must allocate more memory.
755     */
756     size_t total_size = free_list_block_sizes[min_index + 1] + METADATA_TOTAL_AMOUNT;
757     // round up to the size of the next block list size, so we guarantee we end up in the next bucket, to avoid
758     // an issue where the user requests, e.g. 2.1KB, and we allocate 3.9KB, which means due to the two lookup types,
759     // it gets it in the wrong bucket
760     struct block* sys_block = RequestBlock(total_size, flags);
761
762     /*
763     * Put the new memory in the free list (which ought to be empty, as wouldn't need to
764     * request new memory otherwise). Then we can allocate the block.
765     */
766     int sys_index = GetInsertionIndex(GetBlockSize(sys_block) - METADATA_TOTAL_AMOUNT);
767
768     assert(head_list[sys_index] == NULL);
769     head_list[sys_index] = sys_block;
770     return AllocateBlock(head_list[sys_index], sys_index, user_requested_size);
771 }
772
773 /**
774  * Allocates memory on the heap. Unless you *really* know what you're doing, you should always
775  * pass HEAP_NO_FAULT. AllocHeap passes this automatically, but this one doesn't (in case you
776  * want to allocate from the pagable pool).
777  */
778 void* AllocHeapEx(size_t size, int flags) {
779     MAX_IRQL(IRQL_SCHEDULER);
780
781     if (flags & HEAP_ALLOW_PAGING) {
782         if (GetIrql() != IRQL_STANDARD) {
783             PanicEx(PANIC_INVALID_IRQL, "cannot AllocHeapEx(HEAP_ALLOW_PAGING) in irql > IRQL_STANDARD");
784         }
785     }
786
787     /*
788     * We cannot allocate zero blocks (as it would be useless, and couldn't be freed.)
789     * Size cannot be negative as a size_t is an unsigned type.
790     */
791     if (size == 0) {
792         return NULL;
793     }
794
795     if (size >= WARNING_LARGE_REQUEST_SIZE && ((flags & HEAP_ALLOW_PAGING) == 0 || (flags & HEAP_NO_FAULT) != 0)) {
796         LogDeveloperWarning("AllocHeapEx called with allocation of size 0x%X. You should consider using MapVirt.\n", size);
797     }
798
799     if (flags == 0) {
800         LogDeveloperWarning("AllocHeapEx called with flags = 0. You probably meant to pass either HEAP_ALLOW_PAGING,"
801                             "or HEAP_NO_FAULT. Passing neither is valid and it puts it on the locked heap, but allocation"
802                             "may cause faults. This is unlikely to be what you want.");
803     }
804
805     if (flags & HEAP_ALLOW_PAGING) {
806         if (!IsHeapReinitialised()) {
807             PanicEx(PANIC_ASSERTION_FAILURE, "can't allow paging without heap reinit");
808         }
809     }
810
811     bool acquired = LockHeap(flags & HEAP_ALLOW_PAGING);
812
813     size = RoundUpSize(size);
814
815     struct block* block = FindBlock(size, flags);
816
817 #ifndef NDEBUB

```

```

818     outstanding_allocations++;
819 #endif
820
821     if (acquired) {
822         UnlockHeap(flags & HEAP_ALLOW_PAGING);
823     }
824
825     assert(((size_t) block & (ALIGNMENT - 1)) == 0);
826
827     void* ptr = AddVoidPtr(block, METADATA_LEADING_AMOUNT);
828     if (flags & HEAP_ZERO) {
829         inline_memset(ptr, 0, size);
830     }
831
832     return ptr;
833 }
834
835 void* AllocHeap(size_t size) {
836     return AllocHeapEx(size, HEAP_NO_FAULT);
837 }
838
839 void* AllocHeapZero(size_t size) {
840     return AllocHeapEx(size, HEAP_ZERO | HEAP_NO_FAULT);
841 }
842
843 void FreeHeap(void* ptr) {
844     MAX_IQRL(IQRL_SCHEDULER);
845
846     if (ptr == NULL) {
847         return;
848     }
849
850     struct block* block = SubVoidPtr(ptr, METADATA_LEADING_AMOUNT);
851     bool pagable = IsOnSwappableHeap(block);
852     block->prev = NULL;
853     block->next = NULL;
854
855     bool acquired = LockHeap(pagable);
856
857     AddBlock(block);
858
859 #ifndef NDEBUG
860     outstanding_allocations--;
861 #endif
862
863     if (acquired) {
864         UnlockHeap(pagable);
865     }
866 }
867
868 void ReinitHeap(void) {
869     heap_lock = CreateMutex("heap");
870 }
871
872 void InitHeap(void) {
873     InitSpinlock(&heap_spinlock, "heapspin", IQRL_SCHEDULER);
874     InitSpinlock(&heap_locker_lock, "locker", IQRL_HIGH);
875 }

```

File: ./mem/physical.c

```

1
2  /*
3   * mem/physical.c - Physical Memory Manager
4   *
5   * There are two allocation systems in use here. The first is a bitmap system, which has a
6   * bit for each page, which when set, indicates that a page is free. This can be used before
7   * virtual memory is available, and provides O(n) allocation time. After virtual memory is
8   * available, a stack-based system is used to provide O(1) allocation time. The bitmap is still
9   * kept in sync with the stack, to allow detection of double-free conditions and for other uses.
10  *
11  * When physical memory is low, we evict pages before we reach the out of memory condition. This
12  * allows eviction code to allocate physical memory without running out of memory.
13  */
14
15 #include <arch.h>
16 #include <physical.h>
17 #include <diskcache.h>
18 #include <common.h>
19 #include <spinlock.h>
20 #include <assert.h>
21 #include <string.h>
22 #include <iqrl.h>
23 #include <log.h>
24 #include <virtual.h>
25 #include <panic.h>
26
27 /*
28  * How many bits are in each entry of the bitmap allocation array. Should be the number
29  * of bits in a size_t.
30  */
31 #define BITS_PER_ENTRY (sizeof(size_t) * 8)
32
33 /*
34  * The maximum physical memory address we can use, in kilobytes.
35  */
36 #define MAX_MEMORY_KBS ARCH_MAX_RAM_KBS
37
38 /*
39  * The number of pages required to reach the maximum physical memory address.
40  */
41 #define MAX_MEMORY_PAGES (MAX_MEMORY_KBS / ARCH_PAGE_SIZE * 1024)
42
43 /*
44  * The number of entries in the bitmap allocation table required to keep track of
45  * any page up to the maximum usable physical memory address.
46  */
47 #define BITMAP_ENTRIES (MAX_MEMORY_PAGES / BITS_PER_ENTRY)
48
49 /*
50  * We will start evicting pages once we have fewer than this many pages left on the system.
51  * We can have less than this available, if in the process of eviction it causes more pages
52  * to be allocated (this is why we set it to something higher than 0 or 1, to provide a buffer
53  * for eviction to work in).
54  */
55 #define NUM_EMERGENCY_PAGES 32
56
57 /*
58  * One bit per page. Lower bits refer to lower pages. A clear bit indicates
59  * the page is unavailable (allocated / non-RAM), and a set bit indicates the
60  * page is free.
61  */
62 static size_t allocation_bitmap[BITMAP_ENTRIES];
63
64 /*

```

```

65  * Stores pages that are available for us to allocate. If set to NULL, then we have yet to
66  * reinitialise the physical memory manager, and so it cannot be used. The stack grows
67  * upward, and the pointer is incremented after writing the value on a push. The stack stores
68  * physical page numbers (indexes) instead of addresses.
69  */
70  static size_t* allocation_stack = NULL;
71
72  /*
73  * The index in to the allocation_stack bitmap where the next push operation will put
74  * the value.
75  */
76  static size_t allocation_stack_pointer = 0;
77
78  /*
79  * The number of physical pages available (free) remaining in the system. Gets adjusted
80  * on each allocation or deallocation. Gets set during InitPhys() when scanning the system's
81  * memory map.
82  */
83  static size_t pages_left = 0;
84
85  /*
86  * The total number of allocatable pages on the system. Gets set during InitPhys() and ReinitPhys()
87  */
88  static size_t total_pages = 0;
89
90  /*
91  * The highest physical page number that exists on this system. Gets set during InitPhys() when
92  * scanning the system's memory map.
93  */
94  static size_t highest_valid_page_index = 0;
95
96  /*
97  * A lock to prevent concurrent access to the physical memory manager.
98  */
99  static struct spinlock phys_lock;
100
101  static inline bool IsBitmapEntryFree(size_t index) {
102      size_t base = index / BITS_PER_ENTRY;
103      size_t offset = index % BITS_PER_ENTRY;
104
105      return allocation_bitmap[base] & (1 << offset);
106  }
107
108  static inline void AllocateBitmapEntry(size_t index) {
109      size_t base = index / BITS_PER_ENTRY;
110      size_t offset = index % BITS_PER_ENTRY;
111
112      assert(IsBitmapEntryFree(index));
113
114      allocation_bitmap[base] |= (1 << offset);
115  }
116
117  static inline void DeallocateBitmapEntry(size_t index) {
118      size_t base = index / BITS_PER_ENTRY;
119      size_t offset = index % BITS_PER_ENTRY;
120
121      assert(!IsBitmapEntryFree(index));
122
123      allocation_bitmap[base] |= 1 << offset;
124  }
125
126  static inline void PushIndex(size_t index) {
127      assert(index <= highest_valid_page_index);
128      allocation_stack[allocation_stack_pointer++] = index;
129  }
130
131  static inline size_t PopIndex(void) {
132      assert(allocation_stack_pointer != 0);
133      return allocation_stack[--allocation_stack_pointer];
134  }
135
136  /*
137  * Removes an entry from the stack by value. Only to be used when absolutely required,
138  * as it has O(n) runtime and is therefore very slow.
139  */
140  static void RemoveStackEntry(size_t index) {
141      for (size_t i = 0; i < allocation_stack_pointer; ++i) {
142          if (allocation_stack[i] == index) {
143              memmove(allocation_stack + i, allocation_stack + i + 1, (--allocation_stack_pointer - i) * sizeof(size_t));
144              return;
145          }
146      }
147  }
148
149  /**
150  * Deallocates a page of physical memory that was allocated with AllocPhys(). Does not affect virtual mappings -
151  * that should be taken care of before deallocating.
152  *
153  * @param addr The address of the page to deallocate. Must be page-aligned.
154  */
155  void DeallocPhys(size_t addr) {
156      MAX_IRQL(IRQL_SCHEDULER);
157      assert(addr % ARCH_PAGE_SIZE == 0);
158
159      size_t page = addr / ARCH_PAGE_SIZE;
160
161      AcquireSpinlockIrql(&phys_lock);
162
163      --pages_left;
164      DeallocateBitmapEntry(page);
165      if (allocation_stack != NULL) {
166          PushIndex(page);
167      }
168
169      if (pages_left > NUM_EMERGENCY_PAGES * 2) {
170          SetDiskCaches(DISKCACHE_NORMAL);
171      }
172
173      ReleaseSpinlockIrql(&phys_lock);
174  }
175
176  /**
177  * Deallocates a section of physical memory that was allocated with AllocPhysContinuous(). The entire block
178  * of memory must be deallocated at once, i.e. the start address of the memory should be passed in. Does not
179  * affect virtual mappings - that should be taken care of before deallocating.
180  *
181  * @param addr The address of the section of memory to deallocate. Must be page-aligned.
182  * @param size The size of the allocation. This should be the same value that was passed into AllocPhysContinuous().
183  */
184  void DeallocPhysContiguous(size_t addr, size_t bytes) {
185      MAX_IRQL(IRQL_SCHEDULER);
186
187      size_t pages = BytesToPages(bytes);
188      for (size_t i = 0; i < pages; ++i) {
189          DeallocPhys(addr);
190          addr += ARCH_PAGE_SIZE;
191      }
192  }
193
194  static void EvictPagesIfNeeded(void* context) {

```

```

195 (void) context;
196
197 EXACT_IRQL(IRQL_STANDARD);
198
199 /*
200 * We can't fault later on, so we evict now if we are getting low on memory. If this faults,
201 * the recursion will not cause the spinlock to be re-acquired, and so the evicted code won't
202 * run again either - this prevents infinite recursion loops. These fault handlers and recursive
203 * calls can allocate and make use of these 'emergency' pages that we keep by doing this eviction
204 * before we actually run out of memory.
205 *
206 * We loop so that if the first evictions end up needing to allocate memory, we can hopefully
207 * perform another eviction to make up for it (that shouldn't need extra memory).
208 */
209
210 // TODO: probs needs lock on pages_left
211
212 extern int handling_page_fault;
213 if (handling_page_fault > 0) {
214     return;
215 }
216
217 if (pages_left < NUM_EMERGENCY_PAGES) {
218     SetDiskCaches(DISKCACHE_TOSS);
219
220 } else if (pages_left < NUM_EMERGENCY_PAGES * 3 / 2) {
221     SetDiskCaches(DISKCACHE_REDUCE);
222 }
223
224 int timeout = 0;
225 while (pages_left < NUM_EMERGENCY_PAGES && timeout < 5) {
226     handling_page_fault++;
227     EvictVirt();
228     handling_page_fault--;
229     ++timeout;
230 }
231
232 if (pages_left == 0) {
233     Panic(PANIC_OUT_OF_PHYS);
234 }
235
236
237 /**
238 * Allocates a page of physical memory. The resulting memory can be freed with DeallocPhys(). May cause
239 * pages to be evicted from RAM in order to have enough physical memory. Direct users of this function
240 * should be careful - any physical pages that have been allocated but not mapped into virtual memory
241 * cannot be swapped out, and therefore they should be mapped into virtual memory.
242 *
243 * @return The start address of the page of physical memory, or 0 if a page could not be allocated.
244 */
245 size_t AllocPhys(void) {
246     MAX_IRQL(IRQL_SCHEDULER);
247
248     AcquireSpinlockIrql(&phys_lock);
249
250     if (pages_left <= NUM_EMERGENCY_PAGES) {
251         LogWriteSerial("deferring EvictPagesIfNeeded (pages_left = %d)\n", pages_left);
252         DeferUntilIrql(IRQL_STANDARD, EvictPagesIfNeeded, NULL);
253     }
254
255     if (pages_left == 0) {
256         Panic(PANIC_OUT_OF_PHYS);
257     }
258
259     size_t index = 0;
260     if (allocation_stack == NULL) {
261         /*
262          * No stack yet, so must use the bitmap. We could optimise and keep track of the most recently
263          * returned index, but this code is only used during InitVirt(), so we'll try to keep the code
264          * here short and simple.
265          */
266         while (!IsBitmapEntryFree(index)) {
267             index = (index + 1) % MAX_MEMORY_PAGES;
268         }
269     } else {
270         index = PopIndex();
271     }
272
273     AllocateBitmapEntry(index);
274     --pages_left;
275
276     if (pages_left == 0) {
277         LogDeveloperWarning("THAT WAS THE LAST PAGE!\n");
278     }
279
280     ReleaseSpinlockIrql(&phys_lock);
281
282     return index * ARCH_PAGE_SIZE;
283 }
284
285 /**
286 * Allocates a section of contiguous physical memory, that may or may not have requirements as
287 * to where the memory can be located. Allocation in this way is very slow, and so should only
288 * be called where absolutely necessary (e.g. initialising drivers). Must only be called after
289 * a call to ReinitPhys() is made. Deallocation should be done by DeallocPhysContiguous(), passing
290 * in the same size value as passed into AllocPhysContiguous() on allocation. Will not cause pages
291 * to be evicted from RAM, so sufficient memory must exist on the system for this allocation to
292 * succeed.
293 *
294 * @param bytes The size of the allocation, in bytes.
295 * @param min_addr The allocated memory region will not contain any addresses that are lower than
296 * this value.
297 * @param max_addr The allocated memory region will not contain any addresses that are greater than
298 * or equal to this value. If there is no maximum, set this to 0.
299 * @param boundary The allocated memory will not contain any addresses that are an integer multiple
300 * of this value (although it may start at an integer multiple of this address).
301 * If there are no boundary requirements, set this to 0.
302 * @return The start address of the returned physical memory area. If the request could not be
303 * satisfied (e.g. out of memory, no contiguous block, cannot meet requirements), then 0
304 * is returned.
305 */
306 size_t AllocPhysContiguous(size_t bytes, size_t min_addr, size_t max_addr, size_t boundary) {
307     /*
308      * This function should not be called before the stack allocator is setup.
309      * (There is no need for InitVirt() to use this function, and so checking here removes
310      * a check that would have to be done in a loop later).
311      */
312     if (allocation_stack == NULL) {
313         return 0;
314     }
315
316     size_t pages = BytesToPages(bytes);
317     size_t min_index = (min_addr + ARCH_PAGE_SIZE - 1) / ARCH_PAGE_SIZE;
318     size_t max_index = max_addr == 0 ? highest_valid_page_index + 1 : max_addr / ARCH_PAGE_SIZE;
319     size_t count = 0;
320
321     AcquireSpinlockIrql(&phys_lock);
322
323     /*
324      * We need to check we won't try to over-allocate memory, or allocate so much memory that it puts

```

```

325     * us in a critical position.
326     */
327     if (pages + NUM_EMERGENCY_PAGES >= pages_left) {
328         ReleaseSpinlockIrql(&phys_lock);
329         return 0;
330     }
331 }
332 for (size_t index = min_index; index < max_index; ++index) {
333     /*
334      * Reset the counter if we are no longer contiguous, or if we have cross a boundary
335      * that we can't cross.
336      */
337     if (!IsBitmapEntryFree(index) || (boundary != 0 && (index % (boundary / ARCH_PAGE_SIZE) == 0))) {
338         count = 0;
339         continue;
340     }
341     ++count;
342     if (count == pages) {
343         /*
344          * Go back to the start of the section and mark it all as allocated.
345          */
346         size_t start_index = index - count + 1;
347         while (start_index <= index) {
348             AllocateBitmapEntry(start_index);
349             RemoveStackEntry(start_index);
350             ++start_index;
351         }
352         ReleaseSpinlockIrql(&phys_lock);
353         return start_index * ARCH_PAGE_SIZE;
354     }
355 }
356 }
357
358 ReleaseSpinlockIrql(&phys_lock);
359 return 0;
360 }
361
362 /**
363  * Initialises the physical memory manager for the first time. Must be called before any other
364  * memory management function is called. It determines what memory is available on the system
365  * and prepares the O(n) bitmap allocator. This will be slow, but is only needed until ReinitHeap()
366  * gets called. Must only be called once.
367  */
368 void InitPhys(void) {
369     InitSpinlock(&phys_lock, "phys", IRQL_SCHEDULER);
370 }
371
372 /*
373  * Scan the memory tables and fill in the memory that is there.
374  */
375 while (true) {
376     struct arch_memory_range* range = ArchGetMemory();
377     if (range == NULL) {
378         /* No more memory exists */
379         break;
380     } else {
381         /*
382          * Round conservatively (i.e., round the first page up, and the last page down)
383          * so we don't accidentally allow non-existent memory to be allocated.
384          */
385         size_t first_page = (range->start + ARCH_PAGE_SIZE - 1) / ARCH_PAGE_SIZE;
386         size_t last_page = (range->start + range->length) / ARCH_PAGE_SIZE;
387         while (first_page < last_page && first_page < MAX_MEMORY_PAGES) {
388             DeallocateBitmapEntry(first_page);
389             ++pages_left;
390             ++total_pages;
391             if (first_page > highest_valid_page_index) {
392                 highest_valid_page_index = first_page;
393             }
394             ++first_page;
395         }
396     }
397 }
398
399 static void ReclaimBitmapSpace(void) {
400     /*
401      * We can save a tiny bit of extra physical memory on low-memory systems by deallocating the memory
402      * in the bitmap that can't be reached (due to the system not having memory that goes up that high).
403      * e.g. if we allow the system to access 16GB of RAM, but we only have 4MB, then we can save 31 pages
404      * (or 124KB), which on a system with only 4MB of RAM is 3% of total physical memory).
405      * Be careful! If we do this incorrectly we will get memory corruption and some *very* mysterious bugs.
406      */
407     size_t num_unreachable_pages = MAX_MEMORY_PAGES - (highest_valid_page_index + 1);
408     size_t num_unreachable_entries = num_unreachable_pages / BITS_PER_ENTRY;
409     size_t num_unreachable_bitmap_pages = num_unreachable_entries / ARCH_PAGE_SIZE;
410     size_t end_bitmap = ((size_t) allocation_bitmap) + sizeof(allocation_bitmap);
411     /*
412      * DO NOT ROUND UP, or variables in the same page as the end of the bitmap will also be counted as 'free',
413      * causing kernel memory corruption for whatever comes in RAM after that.
414      */
415     size_t unreachable_region = (end_bitmap - ARCH_PAGE_SIZE * num_unreachable_bitmap_pages) % ~(ARCH_PAGE_SIZE - 1);
416     while (num_unreachable_bitmap_pages-- > 0) {
417         DeallocPhys(ArchVirtualToPhysical(unreachable_region));
418         unreachable_region += ARCH_PAGE_SIZE;
419         ++total_pages;
420     }
421 }
422
423 /**
424  * Reinitialises the physical memory manager with a constant-time page allocation system.
425  * Must be called after virtual memory has been initialised. Must only be called once. Must
426  * be called before calling AllocPageContiguous() is called. Should be called as soon as
427  * possible after virtual memory is available.
428  */
429 void ReinitPhys(void) {
430     assert(allocation_stack == NULL);
431     allocation_stack = (size_t*) MapVirt(0, 0, (highest_valid_page_index + 1) * sizeof(size_t), VM_READ | VM_WRITE | VM_LOCK, NULL, 0);
432     for (size_t i = 0; i < MAX_MEMORY_PAGES; ++i) {
433         if (IsBitmapEntryFree(i)) {
434             PushIndex(i);
435         }
436     }
437     ReclaimBitmapSpace();
438 }
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454

```

```

455 size_t GetTotalPhysKilobytes(void) {
456     return total_pages * (ARCH_PAGE_SIZE / 1024);
457 }
458
459 size_t GetFreePhysKilobytes(void) {
460     return pages_left * (ARCH_PAGE_SIZE / 1024);
461 }

```

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     return COMPARE_SIGN(a_entry->virtual, b_entry->virtual);
88 }
89
90 /**
91  * Initialises a virtual address space in an already allocated section of memory.
92  *
93  * @param vas The memory to initialise a virtual address space object in.
94  * @param flags Can be 0 or VAS_NO_ARCH_INIT. If VAS_NO_ARCH_INIT is provided, then no architecture-specific
95  * code will be called. This flag should only be set if called by architecture-specific functions
96  * (e.g. to create the initial address space). Normally, 0 should be passed in.
97  *
98  * @maxirq1 IRQ1_SCHEDULER
99  */
100 void CreateVasEx(struct vas* vas, int flags) {
101     MAX_IRQ1_IRQ1_SCHEDULER;
102
103     vas->mappings = AvlTreeCreate();
104
105     /*
106      * We are in for a world of hurt if someone is able to page fault while
107      * holding the lock on a virtual address space, so better make it IRQ1_SCHEDULER.
108      */
109     InitSpinlock(&vas->lock, "vas", IRQ1_SCHEDULER);
110     AvlTreeSetComparator(vas->mappings, VirtAvlComparator);
111     if (!(flags & VAS_NO_ARCH_INIT)) {
112         ArchInitVas(vas);
113     }
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  * @maxirq1 IRQ1_SCHEDULER
123  */
124 struct vas* CreateVas() {
125     MAX_IRQ1(IRQ1_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     /*
218     * Don't perform relocations on the first load, as the first load will be when 'proper' relocation
219     * happens (i.e. the 'all at once' relocations) - and therefore the quick relocation table will not
220     * be created yet and we'll crash.
221     *
222     * The reason we can't just not do the initial big relocation and make it all work though demand loading
223     * is because not all pages with driver code/data end up being marked as VM_RELOCATABLE (e.g. for small
224     * parts of data segments, etc.).
225     */
226     bool needs_relocations = entry->relocatable && !entry->first_load;
227
228     ArchUpdateMapping(vas, entry);
229     ArchFlushTlb(vas);
230
231     /*
232     * Need to keep page locked if we're doing relocations on it - otherwise by the time that we actually
233     * load in all the data we need to do the relocations (e.g. ELF headers, the symbol table), we have probably
234     * already swapped out the page we are relocating (which leads to us getting nowhere).
235     */
236     if (!needs_relocations) {
237         entry->first_load = false;
238         entry->load_in_progress = false;
239         entry->lock = false;
240     }
241     ReleaseSpinlockIrql(&vas->lock);
242
243     UnmapVirt(target_address, ARCH_PAGE_SIZE);
244
245     if (needs_relocations) {
246         LogWriteSerial(" ----> ABOUT TO PERFORM RELOCATION FIXUPS\n");
247     }

```

```

247     PerformDriverRelocationOnPage(vas, entry->relocation_base, access->address);
248     AcquireSpinlockIrql(&vas->lock);
249     entry->first_load = false;
250     entry->load_in_progress = false;
251     UnlockVirtEx(vas, access->address);
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         if (entry->write && entry->relocatable) {
316             DeferDiskWrite(entry->virtual, entry->file_node, entry->file_offset);
317         }
318
319         entry->in_ram = false;
320         entry->allocated = false;
321         DeallocPhys(entry->physical);
322         ArchUnmap(vas, entry);
323         ArchFlushTlb(vas);
324     } else {
325         /*
326          * Otherwise, we need to mark it as swapfile.
327          */
328         entry->in_ram = false;
329         entry->swapfile = true;
330         entry->allocated = false;
331
332         uint64_t offset = AllocateSwapfileIndex() * ARCH_PAGE_SIZE;
333
334         //PutsConsole("PAGE OUT\n");
335         LogWriteSerial(" ----> WRITING VIRT 0x%X TO SWAP: DISK INDEX 0x%X (offset 0x%X)\n", entry->virtual, (int) offset / ARCH_PAGE_SIZE, (int) offset);
336         DeferDiskWrite(entry->virtual, GetSwapfile(), offset);
337         entry->swapfile_offset = offset;
338
339         ArchUnmap(vas, entry);
340         DeallocPhys(entry->physical);
341         ArchFlushTlb(vas);
342     }
343
344     ReleaseSpinlockIrql(&vas->lock);
345 }
346
347 /**
348  * Lower value means it should be swapped out first.
349  */
350 static int GetPageEvictionRank(struct vas* vas, struct vas_entry* entry) {
351     (void) vas;
352
353     /*
354      * Want to evict in this order:
355      * - file and non-writable
356      * - file and writable
357      * - non-writable
358      * - writable
359      *
360      * When we have a way of dealing with accessed / dirty, it should be in this order:
361      *
362      * 0 VM EVICT FIRST
363      * 10 FILE, NON-WRITABLE, NON-ACCESSED
364      * 20 FILE, WRITABLE, NON-DIRTY, NON-ACCESSED
365      * 30 FILE, NON-WRITABLE, ACCESSED
366      * 40 FILE, WRITABLE, NON-DIRTY, ACCESSED
367      * 50 NORMAL, NON-DIRTY, NON-ACCESSED
368      * 60 NORMAL, NON-DIRTY, ACCESSED
369      * 70 FILE, WRITABLE, DIRTY
370      * 80 NORMAL, DIRTY
371      * 90 COW
372      * 150 RELOCATABLE
373     */
374 }

```



```

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 /**
486  * Searches through virtual memory (that doesn't necessarily have to be in the current virtual address space),
487  * and finds and evicts a page of virtual memory, to try free up physical memory.
488  *
489  * @maxirq1 IRQL_STANDARD
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
540 if (lowest_ranked.entry != NULL) {
541     previous_swaps[swap_num++] = &lowest_ranked.entry;
542     EvictPage(lowest_ranked.vas, lowest_ranked.entry);
543     lowest_ranked.entry->times_swapped++;
544 }
545 }
546
547 static void InsertIntoAvl(struct vas* vas, struct vas_entry* entry) {
548     assert(IsSpinlockHeld(&vas->lock));
549
550     if (entry->global) {
551         AcquireSpinlockIrql(&GetCpu()->global_mappings_lock);
552         AvlTreeInsert(GetCpu()->global_vas_mappings, entry);
553         ReleaseSpinlockIrql(&GetCpu()->global_mappings_lock);
554     }
555     else {
556         AvlTreeInsert(vas->mappings, entry);
557     }
558 }
559
560 static void DeleteFromAvl(struct vas* vas, struct vas_entry* entry) {
561     assert(IsSpinlockHeld(&vas->lock));
562     if (entry->global) {
563         AcquireSpinlockIrql(&GetCpu()->global_mappings_lock);
564         AvlTreeDelete(GetCpu()->global_vas_mappings, entry);
565         ReleaseSpinlockIrql(&GetCpu()->global_mappings_lock);
566     }
567     else {
568         AvlTreeDelete(vas->mappings, entry);
569     }
570 }
571
572 /**
573 * Adds a virtual page mapping to the specified virtual address space. This will add it both to the mapping tree
574 * and the architectural paging structures (so that page faults can be raised, etc., if there is no backing yet).
575 *
576 * @param vas The virtual address space to map this page to
577 * @param physical Only used if VM_LOCK is specified in flags. Determines the physical page that will back the
578 * virtual mapping. If VM_LOCK is set, and this is 0, then a physical page will be allocated. If
579 * VM_LOCK is set, and this is non-zero, then that physical address will be used.
580 * @param virtual The virtual address to map the memory to. This should be non-zero.
581 * @param flags Various bitflags to affect the attributes of the mapping. Flags that are used here are:
582 * VM_READ : if set, the page will be marked as readable
583 * VM_WRITE : if set, the page will be marked as writable
584 * VM_USER : if set, then usermode can access this page without faulting
585 * VM_EXEC : if set, then code can be executed in this page
586 * VM_LOCK : if set, the page will immediately get a physical memory backing, and will not be
587 * paged out
588 * VM_FILE : if set, this page is file-backed
589 * All other flags are ignored by this function.
590 * @param file If VM_FILE is set, then the page is backed by this file, starting at the position specified by pos.
591 * @param pos If VM_FILE is set, then this is the offset into the file where the page is mapped to.
592 *
593 * @maxirql IRQL_SCHEDULER
594 */
595 static void AddMapping(struct vas* vas, size_t physical, size_t virtual, int flags, struct open_file* file, off_t pos, size_t number) {
596     MAX_IRQL(IRQL_SCHEDULER);
597
598     assert(!((file != NULL && (flags & VM_FILE) == 0)));
599
600     struct vas_entry* entry = AllocHeapZero(sizeof(struct vas_entry));
601     entry->allocated = false;
602
603     bool lock = flags & VM_LOCK;
604     entry->lock = lock;
605     entry->in_ram = lock;
606
607     size_t relocation_base = 0;
608     if (flags & VM_RELOCATABLE) {
609         relocation_base = physical;
610         physical = 0;
611     }
612
613     if (lock) {
614         /*
615          * We are not allowed to check if the physical page is allocated/free, because it might come
616          * from a VM_MAP_HARDWARE request, which can map non-RAM pages.
617          */
618         if (physical == 0) {
619             physical = AllocPhys();
620             entry->allocated = true;
621         }
622     }
623
624     /*
625     * MapVirt checks for conflicting flags and returns, so this code doesn't need to worry about that.
626     */
627     entry->virtual = virtual;
628     entry->times_swapped = 0;
629     entry->read = (flags & VM_READ) ? 1 : 0;
630     entry->write = (flags & VM_WRITE) ? 1 : 0;
631     entry->exec = (flags & VM_EXEC) ? 1 : 0;
632     entry->file = (flags & VM_FILE) ? 1 : 0;
633     entry->user = (flags & VM_USER) ? 1 : 0;
634     entry->evict_first = (flags & VM_EVICT_FIRST) ? 1 : 0;
635     entry->relocatable = (flags & VM_RELOCATABLE) ? 1 : 0;
636     entry->allow_temp_write = false;

```

```

637 entry->load_in_progress = false;
638 entry->global = !(flags & VM_LOCAL);
639 entry->physical = physical;
640 entry->ref_count = 1;
641 entry->file_offset = pos;
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, rxwgu'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 InsertIntoAvl(vas, entry);
667 ArchAddMapping(vas, entry);
668
669 if (entry->lock && (flags & VM_MAP_HARDWARE) == 0) {
670     if (GetVas() == vas) {
671         /*
672          * Need to zero out the page - this must happen on first load in, and as we have to load in
673          * locked pages now, we must do it now.
674          */
675         memset((void*) entry->virtual, 0, entry->num_pages * ARCH_PAGE_SIZE);
676     } else {
677         LogDeveloperWarning("yuck. PAGE HAS NOT BEEN ZEROED!\n");
678     }
679 }
680
681 if ((flags & VM_RECURSIVE) == 0) {
682     ReleaseSpinlockIrql(&vas->lock);
683 }
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;
690     dummy.num_pages = 1;
691     dummy.virtual = virtual;
692
693     /*
694      * We have to loop over the local one, and if it isn't there, the global one. We do this
695      * in separate loops to prevent the need to acquire both spinlocks at once, which could lead
696      * to a deadlock.
697      */
698
699     AcquireSpinlockIrql(&vas->lock);
700     for (size_t i = 0; i < pages; ++i) {
701         if (AvlTreeContains(vas->mappings, (void*) &dummy)) {
702             in_use = true;
703             break;
704         }
705         dummy.virtual += ARCH_PAGE_SIZE;
706     }
707     ReleaseSpinlockIrql(&vas->lock);
708
709     if (in_use) {
710         return true;
711     }
712
713     AcquireSpinlockIrql(&GetCpu()->global_mappings_lock);
714     dummy.virtual = virtual;
715     for (size_t i = 0; i < pages; ++i) {
716         if (AvlTreeContains(GetCpu()->global_vas_mappings, (void*) &dummy)) {
717             in_use = true;
718             break;
719         }
720         dummy.virtual += ARCH_PAGE_SIZE;
721     }
722     ReleaseSpinlockIrql(&GetCpu()->global_mappings_lock);
723
724     return in_use;
725 }
726
727 static size_t AllocVirtRange(struct vas* vas, size_t pages, int flags) {
728     /*
729      * TODO: make this deallocatable, and not x86 specific (with that memory address)
730      */
731     if (flags & VM_LOCAL) {
732         /*
733          * Also needs to use the vas to work out what's allocated in that vas
734          */
735         (void) vas;
736         static size_t hideous_allocator = 0x20000000U;
737         size_t retv = hideous_allocator;
738         hideous_allocator += pages * ARCH_PAGE_SIZE;
739         return retv;
740     } else {
741         /*
742          * TODO: this probably needs a global lock of some sort.
743          */
744         static size_t hideous_allocator = ARCH_KRNL_SBRK_BASE;
745         size_t retv = hideous_allocator;
746         hideous_allocator += pages * ARCH_PAGE_SIZE;
747         return retv;
748     }
749 }
750
751 static void FreeVirtRange(struct vas* vas, size_t virtual, size_t pages) {
752     (void) virtual;
753     (void) vas;
754     (void) pages;
755 }
756
757 /**
758  * Creates a virtual memory mapping.
759  *
760  * All mapped pages will be zeroed out (either on first use, or if locked, when allocated) - except if VM_MAP_HARDWARE or
761  * VM_FILE is set. If VM_FILE is set, reading beyond the end of the file, but within the page limit, will read zeroes.
762  *
763  * @param vas The virtual address space to map this page to
764  * @param physical Only used if VM_LOCK is specified in flags. Determines the physical page that will back the
765  * virtual mapping. If VM_LOCK is set, and this is 0, then a physical page will be allocated. If

```

```

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

```

```

897 *
898 * @maxirq1 IRQ1_SCHEDULER
899 */
900 size_t MapVirt(size_t physical, size_t virtual, size_t bytes, int flags, struct open_file* file, off_t pos) {
901     MAX_IRQ1_IRQ1_SCHEDULER;
902     size_t pages = BytesToPages(bytes);
903     size_t ret = MapVirtEx(GetVas(), physical, virtual, pages, flags, file, pos);
904     return ret;
905 }
906
907 static struct vas_entry* GetVirtEntry(struct vas* vas, size_t virtual) {
908     struct vas_entry dummy;
909     dummy.num_pages = 1;
910     dummy.virtual = virtual &~ (ARCH_PAGE_SIZE - 1);
911
912     assert(IsSpinlockHeld(&vas->lock));
913
914     struct vas_entry* res = (struct vas_entry*) AvlTreeGet(vas->mappings, (void*) &dummy);
915     if (res == NULL) {
916         AcquireSpinlockIrql(&GetCpu()->global_mappings_lock);
917         res = (struct vas_entry*) AvlTreeGet(GetCpu()->global_vas_mappings, (void*) &dummy);
918         ReleaseSpinlockIrql(&GetCpu()->global_mappings_lock);
919     }
920     return res;
921 }
922
923 size_t GetPhysFromVirt(size_t virtual) {
924     struct vas* vas = GetVas();
925     AcquireSpinlockIrql(&vas->lock);
926     struct vas_entry* entry = GetVirtEntry(vas, virtual);
927     size_t result = entry->physical;
928
929     /*
930      * Handle mappings of more than 1 page at a time by adding the extra offset
931      * from the start of the mapping.
932      */
933     size_t target_page = virtual / ARCH_PAGE_SIZE;
934     size_t entry_page = entry->virtual / ARCH_PAGE_SIZE;
935     if (entry_page < target_page)
936         result += (target_page - entry_page) * ARCH_PAGE_SIZE;
937 }
938
939 ReleaseSpinlockIrql(&vas->lock);
940 return result;
941 }
942
943 static size_t SplitLargePageEntryIntoMultiple(struct vas* vas, size_t virtual, struct vas_entry* entry, int num_to_leave) {
944     if (entry->num_pages == 1) {
945         return 1;
946     }
947
948     if (entry->ref_count != 1) {
949         LogDeveloperWarning("Splitting multi-mapping with ref_count != 1, this hasn't been tested!\n");
950     }
951
952     /*
953      * Although it can't be allocated, it can be in RAM (e.g. for VM_MAP_HARDWARE).
954      */
955     assert(!entry->allocated);
956     assert(!entry->swapfile);
957
958     size_t entry_page = entry->virtual / ARCH_PAGE_SIZE;
959     size_t target_page = virtual / ARCH_PAGE_SIZE;
960
961     /*
962      * Split off anything before this page.
963      */
964     if (entry_page < target_page) {
965         size_t num_beforehand = target_page - entry_page;
966
967         struct vas_entry* pre_entry = AllocHeap(sizeof(struct vas_entry));
968         *pre_entry = *entry;
969
970         pre_entry->num_pages = num_beforehand;
971         entry->num_pages -= num_beforehand;
972         entry->virtual += num_beforehand * ARCH_PAGE_SIZE;
973
974         /*
975          * For multi-mapping for VM_MAP_HARDWARE
976          */
977         if (entry->physical != 0) {
978             entry->physical += num_beforehand * ARCH_PAGE_SIZE;
979         }
980
981         if (entry->file) {
982             entry->file_offset += num_beforehand * ARCH_PAGE_SIZE;
983         }
984
985         InsertIntoAvl(vas, pre_entry);
986     }
987
988     /*
989      * There's now no pages beforehand. Now we need to check if there are any other pages
990      * after this.
991      */
992     if (entry->num_pages > num_to_leave) {
993         struct vas_entry* post_entry = AllocHeap(sizeof(struct vas_entry));
994         *post_entry = *entry;
995
996         post_entry->num_pages -= num_to_leave;
997         entry->num_pages = num_to_leave;
998
999         post_entry->virtual += ARCH_PAGE_SIZE * num_to_leave;
1000
1001         /*
1002          * For multi-mapping for VM_MAP_HARDWARE
1003          */
1004         if (entry->physical != 0) {
1005             post_entry->physical += ARCH_PAGE_SIZE * num_to_leave;
1006         }
1007
1008         if (entry->file) {
1009             post_entry->file_offset += ARCH_PAGE_SIZE * num_to_leave;
1010         }
1011
1012         InsertIntoAvl(vas, post_entry);
1013     }
1014
1015     return entry->num_pages;
1016 }
1017
1018 static void BringIntoMemoryFromCow(struct vas_entry* entry) {
1019     /*
1020      * If someone deallocates a COW page in another process to get the ref
1021      * count back to 1 already, then we just have the page to ourselves again.
1022      */
1023     if (entry->ref_count == 1) {
1024         entry->cow = false;
1025         ArchUpdateMapping(GetVas(), entry);
1026         ArchFlushTlb(GetVas());
1027     }

```

```

1027     return;
1028 }
1029
1030 uint8_t page_data[ARCH_PAGE_SIZE];
1031 inline_memcpy(page_data, (void*) entry->virtual, ARCH_PAGE_SIZE);
1032
1033 entry->ref_count--;
1034
1035 if (entry->ref_count == 1) {
1036     entry->cow = false;
1037 }
1038
1039 struct vas_entry* new_entry = AllocHeap(sizeof(struct vas_entry));
1040 *new_entry = *entry;
1041 new_entry->ref_count = 1;
1042 new_entry->physical = AllocPhys();
1043 new_entry->allocated = true;
1044 DeleteFromAvl(GetVas(), entry);
1045 FreeHeap(entry);
1046 ArchUpdateMapping(GetVas(), entry);
1047 ArchFlushTlb(GetVas());
1048 inline_memcpy((void*) entry->virtual, page_data, ARCH_PAGE_SIZE);
1049 }
1050
1051 static void BringIntoMemoryFromFile(struct vas_entry* entry, size_t faulting_virt) {
1052     // TODO: need to test that you're allowed to read past the end of the file (even into other pages)
1053     // if the size mapped allows it, and just get zeros
1054
1055     SplitLargePageEntryIntoMultiple(GetVas(), faulting_virt, entry, 1);
1056     entry->load_in_progress = true;
1057     ArchUpdateMapping(GetVas(), entry);
1058     ArchFlushTlb(GetVas());
1059     DeferDiskRead(entry->virtual, entry->file_node, entry->file_offset, false);
1060 }
1061
1062 static void BringIntoMemoryFromSwapfile(struct vas_entry* entry) {
1063     assert(!entry->file);
1064
1065     uint64_t offset = entry->swapfile_offset;
1066     entry->load_in_progress = true;
1067     ArchUpdateMapping(GetVas(), entry);
1068     ArchFlushTlb(GetVas());
1069     DeferDiskRead(entry->virtual, GetSwapfile(), offset, true);
1070 }
1071
1072 static void BringInBlankPage(struct vas* vas, struct vas_entry* entry, size_t faulting_virt, int fault_type) {
1073     if ((fault_type & VM_READ) && !entry->read) {
1074         UnhandledFault();
1075     }
1076     if ((fault_type & VM_WRITE) && !entry->write) {
1077         UnhandledFault();
1078     }
1079     if ((fault_type & VM_EXEC) && !entry->exec) {
1080         UnhandledFault();
1081     }
1082
1083     SplitLargePageEntryIntoMultiple(vas, faulting_virt, entry, 1);
1084     assert(entry->num_pages == 1);
1085
1086     entry->physical = AllocPhys();
1087     entry->allocated = true;
1088     entry->in_ram = true;
1089     entry->allow_temp_write = true;
1090     assert(!entry->swapfile);
1091     ArchUpdateMapping(vas, entry);
1092     ArchFlushTlb(vas);
1093
1094     inline_memset((void*) entry->virtual, 0, ARCH_PAGE_SIZE);
1095     entry->allow_temp_write = false;
1096     ArchUpdateMapping(vas, entry);
1097     ArchFlushTlb(vas);
1098 }
1099
1100 static int BringIntoMemory(struct vas* vas, struct vas_entry* entry, bool allow_cow, size_t faulting_virt, int fault_type) {
1101     (void) vas;
1102     assert(IsSpinlockHeld(&vas->lock));
1103
1104     if (entry->cow && allow_cow) {
1105         assert(entry->num_pages == 1);
1106         BringIntoMemoryFromCow(entry);
1107         return 0;
1108     }
1109
1110     if (entry->file && !entry->in_ram) {
1111         BringIntoMemoryFromFile(entry, faulting_virt);
1112         return 0;
1113     }
1114
1115     if (entry->swapfile) {
1116         assert(entry->num_pages == 1);
1117         BringIntoMemoryFromSwapfile(entry);
1118         return 0;
1119     }
1120
1121     if (!entry->in_ram) {
1122         BringInBlankPage(vas, entry, faulting_virt, fault_type);
1123         return 0;
1124     }
1125
1126     return EINVAL;
1127 }
1128
1129 bool LockVirtEx(struct vas* vas, size_t virtual) {
1130     struct vas_entry* entry = GetVirtEntry(vas, virtual);
1131
1132     if (!entry->in_ram) {
1133         SplitLargePageEntryIntoMultiple(vas, virtual, entry, 1);
1134         int res = BringIntoMemory(vas, entry, true, virtual, 0);
1135         if (res != 0) {
1136             Panic(PANIC_CANNOT_LOCK_MEMORY);
1137         }
1138         assert(entry->in_ram);
1139     }
1140
1141     bool old_lock = entry->lock;
1142     entry->lock = true;
1143     return old_lock;
1144 }
1145
1146 void UnlockVirtEx(struct vas* vas, size_t virtual) {
1147     struct vas_entry* entry = GetVirtEntry(vas, virtual);
1148     SplitLargePageEntryIntoMultiple(vas, virtual, entry, 1);
1149     entry->lock = false;
1150 }
1151
1152 bool LockVirt(size_t virtual) {
1153     struct vas* vas = GetVas();
1154     AcquireSpinlockIrql(&vas->lock);
1155     bool res = LockVirtEx(vas, virtual);
1156     ReleaseSpinlockIrql(&vas->lock);

```

```

1157     return res;
1158 }
1159
1160 void UnlockVirt(size_t virtual) {
1161     struct vas* vas = GetVas();
1162     AcquireSpinlockIrql(&vas->lock);
1163     UnlockVirtEx(vas, virtual);
1164     ReleaseSpinlockIrql(&vas->lock);
1165 }
1166
1167 void SetVirtPermissions(size_t virtual, int set, int clear) {
1168     /*
1169     * Only allow these flags to be set / cleared.
1170     */
1171     if ((set | clear) & ~(VM_READ | VM_WRITE | VM_EXEC | VM_USER)) {
1172         assert(false);
1173         return;
1174     }
1175
1176     struct vas* vas = GetVas();
1177     AcquireSpinlockIrql(&vas->lock);
1178
1179     struct vas_entry* entry = GetVirtEntry(vas, virtual);
1180     if (entry == NULL)
1181         PanicEx(PANIC_ASSERTION_FAILURE, "[SetVirtPermissions] got null back for virt entry");
1182 }
1183
1184 SplitLargePageEntryIntoMultiple(vas, virtual, entry, 1);
1185
1186 entry->read = (set & VM_READ) ? true : (clear & VM_READ ? false : entry->read);
1187 entry->write = (set & VM_WRITE) ? true : (clear & VM_WRITE ? false : entry->write);
1188 entry->exec = (set & VM_EXEC) ? true : (clear & VM_EXEC ? false : entry->exec);
1189 entry->user = (set & VM_USER) ? true : (clear & VM_USER ? false : entry->user);
1190
1191 ArchUpdateMapping(vas, entry);
1192 ArchFlushTlb(vas);
1193
1194 ReleaseSpinlockIrql(&vas->lock);
1195 }
1196
1197 int GetVirtPermissions(size_t virtual) {
1198     struct vas* vas = GetVas();
1199     AcquireSpinlockIrql(&vas->lock);
1200     struct vas_entry* entry_ptr = GetVirtEntry(GetVas(), virtual);
1201     if (entry_ptr == NULL) {
1202         ReleaseSpinlockIrql(&vas->lock);
1203         return 0;
1204     }
1205     struct vas_entry entry = *entry_ptr;
1206     ReleaseSpinlockIrql(&vas->lock);
1207
1208     int permissions = 0;
1209     if (entry.read) permissions |= VM_READ;
1210     if (entry.write) permissions |= VM_WRITE;
1211     if (entry.exec) permissions |= VM_EXEC;
1212     if (entry.lock) permissions |= VM_LOCK;
1213     if (entry.file) permissions |= VM_FILE;
1214     if (entry.user) permissions |= VM_USER;
1215     if (entry.global) permissions |= VM_LOCAL;
1216     if (entry.relocatable) permissions |= VM_RELOCATABLE;
1217
1218     return permissions;
1219 }
1220
1221 int UnmapVirtEx(struct vas* vas, size_t virtual, size_t pages) {
1222     bool needs_tlb_flush = false;
1223
1224     for (size_t i = 0; i < pages; ++i) {
1225         struct vas_entry* entry = GetVirtEntry(vas, virtual + i * ARCH_PAGE_SIZE);
1226         if (entry == NULL)
1227             return ENOMEM;
1228     }
1229
1230     SplitLargePageEntryIntoMultiple(vas, virtual, entry, 1); // TODO: multi-pages
1231
1232     assert(entry->ref_count > 0);
1233     entry->ref_count--;
1234
1235     if (entry->ref_count == 0) {
1236         if (entry->file && entry->write && entry->in_ram) {
1237             DeferDiskWrite(entry->virtual, entry->file_node, entry->file_offset);
1238         }
1239         if (entry->in_ram) {
1240             ArchUnmap(vas, entry);
1241             needs_tlb_flush = true;
1242         }
1243         if (entry->swapfile) {
1244             assert(!entry->allocated);
1245             DeallocateSwapfileIndex(entry->physical / ARCH_PAGE_SIZE);
1246         }
1247         if (entry->allocated) {
1248             assert(!entry->swapfile); // can't be on swap, as putting on swap clears allocated bit
1249             DeallocPhys(entry->physical);
1250         }
1251
1252         DeleteFromAvl(vas, entry);
1253         FreeHeap(entry);
1254         FreeVirtRange(vas, virtual + i * ARCH_PAGE_SIZE, entry->num_pages);
1255     }
1256
1257     if (needs_tlb_flush) {
1258         ArchFlushTlb(vas);
1259     }
1260
1261     return 0;
1262 }
1263
1264 int UnmapVirt(size_t virtual, size_t bytes) {
1265     struct vas* vas = GetVas();
1266     AcquireSpinlockIrql(&vas->lock);
1267     int res = UnmapVirtEx(vas, virtual, BytesToPages(bytes));
1268     ReleaseSpinlockIrql(&vas->lock);
1269     return res;
1270 }
1271
1272 static void CopyVasRecursive(struct avl_node* node, struct vas* new_vas) {
1273     if (node == NULL) {
1274         return;
1275     }
1276
1277     CopyVasRecursive(AvlTreeGetLeft(node), new_vas);
1278     CopyVasRecursive(AvlTreeGetRight(node), new_vas);
1279
1280     struct vas_entry* entry = AvlTreeGetData(node);
1281
1282     if (entry->lock) {
1283         /*
1284         * Got to add the new entry right now. We know it must be in memory as it
1285         * is locked.
1286         */
1287     }

```

```

1287 */
1288 assert(entry->in_ram);
1289
1290 if (entry->allocated) {
1291     /*
1292      * Copy the physical page. We do this by copying the data into a buffer,
1293      * putting a new physical page in the existing VAS and then copying the
1294      * data there. Then the original physical page that was there is free to use
1295      * as the copy.
1296      */
1297     uint8_t page_data[ARCH_PAGE_SIZE]; // TODO: MapVirt this ?
1298     inline_memcpy(page_data, (void*) entry->virtual, ARCH_PAGE_SIZE);
1299     size_t new_physical = entry->physical;
1300     entry->physical = AllocPhys();
1301     ArchUpdateMapping(GetVas(), entry);
1302     ArchFlushTlb(GetVas());
1303     inline_memcpy((void*) entry->virtual, page_data, ARCH_PAGE_SIZE);
1304
1305     struct vas_entry* new_entry = AllocHeap(sizeof(struct vas_entry));
1306     *new_entry = *entry;
1307     new_entry->ref_count = 1;
1308     new_entry->physical = new_physical;
1309     new_entry->allocated = true;
1310     AvlTreeInsert(new_vas->mappings, entry); // don't need to insert global - we're copying so it's already in global
1311     ArchAddMapping(new_vas, entry);
1312
1313 } else {
1314     LogWriteSerial("fork() on a hardware-mapped page is not implemented yet");
1315     PanicEx(PANIC_NOT_IMPLEMENTED, "CopyVasRecursive");
1316 }
1317
1318 } else {
1319     /*
1320      * If it's on swap, it's okay to still mark it as COW, as when we reload we will
1321      * try to do the 'copy'-on-write, and then we will reload from swap, and it will
1322      * then reload and then be copied. Alternatively, if it is read, then it gets brought
1323      * back into memory, but as a COW page still.
1324      *
1325      * BSS memory works fine like this too (but will incur another fault when it is used).
1326      *
1327      * At this stage (where shared memory doesn't exist yet), file mapped pages will also
1328      * be COWed. This means there will two copies of the file in memory should they write
1329      * to it. The final process to release memory will ultimately 'win' and have its changes
1330      * perserved to disk (the others will get overwritten).
1331      */
1332     entry->cow = true;
1333     entry->ref_count++;
1334
1335     // again, no need to add to global - it's already there!
1336     AvlTreeInsert(new_vas->mappings, entry);
1337
1338     ArchUpdateMapping(GetVas(), entry);
1339     ArchAddMapping(new_vas, entry);
1340 }
1341 }
1342
1343 struct vas* CopyVas(void) {
1344     struct vas* vas = GetVas();
1345     struct vas* new_vas = CreateVas();
1346
1347     AcquireSpinlockIrql(&vas->lock);
1348     // no need to change global - it's already there!
1349     CopyVasRecursive(AvlTreeGetRootNode(vas->mappings), new_vas);
1350     ArchFlushTlb(vas);
1351     ReleaseSpinlockIrql(&vas->lock);
1352
1353     return new_vas;
1354 }
1355
1356 struct vas* GetVas(void) {
1357     // TODO: cpu probably needs to have a lock object in it called current_vas_lock, which needs to be held whenever
1358     // someone reads or writes to current_vas;
1359     return GetCpu()->current_vas;
1360 }
1361
1362 void SetVas(struct vas* vas) {
1363     GetCpu()->current_vas = vas;
1364     ArchSetVas(vas);
1365 }
1366
1367 struct vas* GetKernelVas(void) {
1368     return kernel_vas;
1369 }
1370
1371 void InitVirt(void) {
1372     // TODO: cpu probably needs to have a lock object in it called current_vas_lock, which needs to be held whenever
1373     // someone reads or writes to current_vas;
1374
1375     assert(!virt_initialised);
1376     GetCpu()->global_vas_mappings = AvlTreeCreate();
1377     AvlTreeSetComparator(GetCpu()->global_vas_mappings, VirtAvlComparator);
1378     ArchInitVirt();
1379
1380     kernel_vas = GetVas();
1381     virt_initialised = true;
1382 }
1383
1384 /**
1385  * Handles a page fault. Only to be called by the low-level, platform specific interrupt handler when a page
1386  * fault occurs. It will attempt to resolve any fault (e.g. handling copy-on-write, swapfile, file-backed, etc.).
1387  *
1388  * @param faulting_virt The virtual address that was accessed that caused the page fault
1389  * @param fault_type The reason why a page fault occurred. Is a bitfield of VM_WRITE, VM_READ, VM_USER and VM_EXEC.
1390  * VM_READ should be set if a non-present page was accessed. VM_USER should be set for permission
1391  * faults, and VM_WRITE should be set if the operation was caused by a write (as opposed to a read).
1392  * VM_EXEC should be set if execution tried to occur in a non-executable page.
1393  *
1394  * @maxirq Irql_PAGE_FAULT
1395  */
1396 int handling_page_fault = 0;
1397
1398 void HandleVirtFault(size_t faulting_virt, int fault_type) {
1399     if (GetIrql() >= IRQL_SCHEDULER) {
1400         PanicEx(PANIC_INVALID_IRQ, "page fault while Irql >= IRQL_SCHEDULER. is some clown holding a spinlock while "
1401             "executing pageable code? or calling AllocHeapEx wrong with a lock held?");
1402     }
1403
1404     struct vas* vas = GetVas();
1405     AcquireSpinlockIrql(&vas->lock);
1406     ++handling_page_fault;
1407
1408     struct vas_entry* entry = GetVirtEntry(vas, faulting_virt);
1409
1410     if (entry == NULL) {
1411         UnhandledFault();
1412     }
1413
1414     if (entry->load_in_progress) {
1415         --handling_page_fault;
1416     }

```



```

1417     ReleaseSpinlockIrql(&vas->lock);
1418     Schedule();
1419     return;
1420 }
1421
1422 /*
1423  * Sanity check that our flags are configured correctly.
1424  */
1425 assert(!(entry->in_ram && entry->swapfile));
1426 assert(!(entry->file && entry->swapfile));
1427 assert(!(entry->in_ram && entry->lock));
1428 assert(!(entry->cow && entry->lock));
1429
1430 // TODO: check for access violations (e.g. user using a supervisor page)
1431 //       (read / write is not necessarily a problem, e.g. COW)
1432
1433 int result = BringIntoMemory(vas, entry, fault_type & VM_WRITE, faulting_virt, fault_type);
1434 if (result != 0) {
1435     UnhandledFault();
1436 }
1437
1438 --handling_page_fault;
1439 ReleaseSpinlockIrql(&vas->lock);
1440 }
1441
1442 /**
1443  * Determines whether or not virtual memory has been initialised yet. This can be used to determine if it
1444  * is possible to call any virtual memory functions (e.g. in the physical memory and heap allocators).
1445  *
1446  * @return True if virtual memory is available, false otherwise.
1447  *
1448  * @maxirql IRQL_HIGH
1449  */
1450 bool IsVirtInitialised(void) {
1451     return virt_initialised;
1452 }
1453
1454 size_t BytesToPages(size_t bytes) {
1455     return (bytes + ARCH_PAGE_SIZE - 1) / ARCH_PAGE_SIZE;
1456 }
1457
1458 void DestroyVas(struct vas* vas) {
1459     /*
1460      * TODO: implement this
1461      */
1462     (void) vas;
1463
1464     if (vas == GetVas()) {
1465         Panic(PANIC_VAS_TRIED_TO_SELF_DESTRUCT);
1466     }
1467
1468     // TODO: may need to add reference counting later on (depending on what we need),
1469     //       just decrement here, and only delete if got to 0.
1470 }

```

File: /debug/DS_Store

[binary]

File: /debug/framework/hostio.c

```

1
2 #ifndef NDEBBUG
3
4 #include <common.h>
5 #include <debug/hostio.h>
6 #include <log.h>
7 #include <assert.h>
8
9 static void outb(uint16_t port, uint8_t value)
10 {
11     asm volatile ("outb %0, %1" : : "a"(value), "Nd"(port));
12 }
13
14 static uint8_t inb(uint16_t port)
15 {
16     uint8_t value;
17     asm volatile ("inb %1, %0"
18 : "=a"(value)
19 : "Nd"(port));
20     return value;
21 }
22
23 static void DbgWriteByte(uint8_t value) {
24     while ((inb(0x3E8 + 5) & 0x20) == 0) {
25     }
26 }
27 outb(0x3E8, value);
28 }
29
30 static uint8_t DbgReadByte() {
31     while ((inb(0x2F8 + 5) & 1) == 0) {
32     }
33 }
34 return inb(0x2F8);
35 }
36 }
37
38 void DbgWritePacket(int type, uint8_t* data, int size) {
39     DbgWriteByte(0xAA);
40     DbgWriteByte(type);
41     DbgWriteByte((size >> 16) & 0xFF);
42     DbgWriteByte((size >> 8) & 0xFF);
43     DbgWriteByte((size >> 0) & 0xFF);
44     DbgWriteByte(0xBB);
45
46     for (int i = 0; i < size; ++i) {
47         DbgWriteByte(data[i]);
48     }
49
50     DbgWriteByte(0xCC);
51 }
52
53 void DbgReadPacket(int* type, uint8_t* data, int* size) {
54     uint8_t v = DbgReadByte();
55     assert(v == 0xAA);
56
57     *type = DbgReadByte();
58
59     int size_tmp = DbgReadByte();
60     size_tmp <<= 8;
61     size_tmp |= DbgReadByte();
62     size_tmp <<= 8;
63     size_tmp |= DbgReadByte();
64
65     int max_size = *size;
66
67     *size = size_tmp;
68
69     v = DbgReadByte();
70     assert(v == 0xBB);
71
72     for (int i = 0; i < size_tmp; ++i) {
73         if (i < max_size) {
74             data[i] = DbgReadByte();
75         }
76     }
77
78     v = DbgReadByte();
79     assert(v == 0xCC);
80 }
81
82 #else
83 extern int make_iso_compilers_happy;
84 #endif

```

File: ./debug/framework/tfw.c

```

1
2 /*
3  * tfw - Testing Framework
4  */
5
6 #include <debug/tfw.h>
7 #include <debug/hostio.h>
8 #include <string.h>
9 #include <log.h>
10 #include <assert.h>
11 #include <panic.h>
12 #include <debug/tfw_tests.h>
13
14 #ifndef NDEBBUG
15
16 #define RESULT_NOT_STARTED 0
17 #define RESULT_IN_PROGRESS 1
18 #define RESULT_SUCCESS 2
19 #define RESULT_FAILURE 3
20 #define RESULT_SKIPPED 4
21
22 struct host_state {
23     int test_num;
24     uint8_t test_results[MAX_TWF_TESTS];
25 };
26
27 static struct host_state test_state;
28 static struct tfw_test registered_tests[MAX_TWF_TESTS];
29 static int num_tests_registered = 0;
30 static bool nightly_mode = false;
31
32 #define PACKET_BUFFER_SIZE (MAX_TWF_TESTS + 256)
33 uint8_t packet_buffer[PACKET_BUFFER_SIZE];
34
35 static bool all_tests_done = false;
36
37 /*
38  * Send a 0x11 byte to host to ask for current data.
39  * Host sends back a packet with 0x22, then the data at data + 8 onwards.
40  * It may also send back 0x55, in which case it means it the start of the run (no data yet).

```

```

41  */
42  static void GetHostState() {
43      if (all_tests_done) {
44          return;
45      }
46
47      packet_buffer[0] = 0x11;
48      DbgWritePacket(DBGPKT_TFW, packet_buffer, 1);
49
50      while (true) {
51          int type;
52          int size = PACKET_BUFFER_SIZE - 32;
53          DbgReadPacket(&type, packet_buffer, &size);
54          if (size < 1024 && type == DBGPKT_TFW) {
55              if (packet_buffer[0] == 0x55) {
56                  test_state.test_num = 0;
57                  inline_memset(test_state.test_results, 0, sizeof(test_state.test_results));
58                  break;
59              }
60              if (packet_buffer[0] == 0x66 || packet_buffer[0] == 0x67) {
61                  nightly_mode = packet_buffer[0] == 0x67;
62                  continue;
63              }
64              assert(packet_buffer[0] == 0x22);
65              test_state = *(struct host_state*) (packet_buffer + 8);
66              break;
67          }
68      }
69  }
70
71  static void ReadAck(void) {
72      int type;
73      int size = 10;
74      uint8_t d[32];
75      DbgReadPacket(&type, d, &size);
76      assert(d[0] == 0x66 || d[0] == 0x67);
77  }
78
79  /*
80   * To set the host state, send a 0x33/0x44, then the data at +8.
81   * 0x33 - starting a test
82   * 0x34 - finished test successfully
83   * 0x35 - finished test unsuccessfully
84   * 0x44 - finished last test successfully
85   * 0x45 - finished last test unsuccessfully
86   */
87  static void SetHostState(int code) {
88      if (all_tests_done) {
89          return;
90      }
91
92      assert(sizeof(struct host_state) < 4096 - 8);
93
94      inline_memset(packet_buffer, 0, sizeof(test_state) + 8 + (code == 0x33 ? MAX_NAME_LENGTH : 0));
95      packet_buffer[0] = code;
96      inline_memcpy(packet_buffer + 8, &test_state, sizeof(test_state));
97
98      if (code == 0x33) {
99          strncpy((char*) (packet_buffer + 8 + sizeof(test_state)), registered_tests[test_state.test_num].name, MAX_NAME_LENGTH);
100      }
101      packet_buffer[1] = registered_tests[test_state.test_num].nightly_only ? 1 : 0;
102
103      DbgWritePacket(DBGPKT_TFW, packet_buffer, sizeof(test_state) + 8 + (code == 0x33 ? MAX_NAME_LENGTH : 0));
104      ReadAck();
105
106      if (code >> 4 == 0x4) {
107          LogWriteSerial("\n\nFINISHED ALL TESTS:\n");
108          for (int i = 0; i < num_tests_registered; ++i) {
109              LogWriteSerial("    %s - %s\n", registered_tests[i].name, test_state.test_results[i] == RESULT_SUCCESS ? "passed" : "failed");
110          }
111          LogWriteSerial("\n");
112          all_tests_done = true;
113      }
114  }
115
116  static bool in_test = false;
117
118  bool IsInTfwTest(void) {
119      return in_test;
120  }
121
122  void RegisterTfwTest(const char* name, int start_point, void (*code)(struct tfw_test*, size_t), int expected_panic, size_t context) {
123      struct tfw_test test;
124      inline_memset(&test, 0, MAX_NAME_LENGTH);
125      strncpy(test.name, name, MAX_NAME_LENGTH - 1);
126      test.code = code;
127      test.expected_panic_code = expected_panic;
128      test.start_point = start_point;
129      test.context = context;
130      test.nightly_only = false;
131      registered_tests[num_tests_registered++] = test;
132  }
133
134  void RegisterNightlyTfwTest(const char* name, int start_point, void (*code)(struct tfw_test*, size_t), int expected_panic, size_t context) {
135      RegisterTfwTest(name, start_point, code, expected_panic, context);
136      registered_tests[num_tests_registered - 1].nightly_only = true;
137  }
138
139  void FinishedTfwTest(int panic_code) {
140      in_test = false;
141
142      bool success = registered_tests[test_state.test_num].expected_panic_code == panic_code;
143      LogWriteSerial("FinishedTfwTest: finished test %d, expected %d vs. actual %d\n", test_state.test_num, registered_tests[test_state.test_num].expected_panic_code, panic_code);
144
145      test_state.test_results[test_state.test_num] = success ? RESULT_SUCCESS : RESULT_FAILURE;
146
147      if (registered_tests[test_state.test_num].nightly_only && !nightly_mode) {
148          test_state.test_results[test_state.test_num] = RESULT_SKIPPED;
149      }
150
151      test_state.test_num++;
152      if (test_state.test_num >= num_tests_registered) {
153          SetHostState(success ? 0x44 : 0x45);
154      } else {
155          SetHostState(success ? 0x34 : 0x35);
156      }
157  }
158
159  void MarkTfwStartPoint(int id) {
160      LogWriteSerial("Reached TFW_SP %d\n", id);
161
162      if (test_state.test_num >= num_tests_registered || all_tests_done) {
163          return;
164      }
165
166      if (registered_tests[test_state.test_num].start_point == id) {
167          test_state.test_results[test_state.test_num] = RESULT_IN_PROGRESS;
168          SetHostState(0x33);
169          in_test = true;
170          if (!registered_tests[test_state.test_num].nightly_only || nightly_mode) {

```

```

171         LogWriteSerial("MarkTfwStartPoint: running test %d\n", test_state.test_num);
172         registered_tests[test_state.test_num].code(registered_tests + test_state.test_num, registered_tests[test_state.test_num].context);
173     }
174     Panic(PANIC_UNIT_TEST_OK);
175 }
176
177
178 void RegisterTfwTests(void) {
179     RegisterTfwInitTests();
180     RegisterTfwWaitTests();
181     RegisterTfwSemaphoreTests();
182     RegisterTfwPhysTests();
183     RegisterTfwAVLTreeTests();
184     RegisterTfwPriorityQueueTests();
185     RegisterTfwIrqTests();
186 }
187
188 void InitTfw(void) {
189     RegisterTfwTests();
190     ReadAck();
191     GetHostState();
192     test_state.test_results[test_state.test_num] = RESULT_NOT_STARTED;
193 }
194
195 #endif

```

File: ./debug/tests/avl.c

```

1
2 #include <debug.h>
3 #include <assert.h>
4 #include <panic.h>
5 #include <string.h>
6 #include <log.h>
7 #include <avl.h>
8 #include <heap.h>
9 #include <physical.h>
10
11 #ifndef NDEBBUG
12
13 /*
14 typedef void (*avl_deletion_handler)(void*);
15 typedef int (*avl_comparator)(void*, void*);
16
17 void* AvlTreeGet(struct avl_tree* tree, void* data);
18 struct avl_node* AvlTreeGetRootNode(struct avl_tree* tree);
19 struct avl_node* AvlTreeGetLeft(struct avl_node* node);
20 struct avl_node* AvlTreeGetRight(struct avl_node* node);
21 void* AvlTreeGetData(struct avl_node* node);
22 avl_deletion_handler AvlTreeSetDeletionHandler(struct avl_tree* tree, avl_deletion_handler handler);
23 avl_comparator AvlTreeSetComparator(struct avl_tree* tree, avl_comparator comparator);*/
24
25 TFW_CREATE_TEST(AVLTreeBasic) { TFW_IGNORE_UNUSED
26     int heap_allocations = DbgGetOutstandingHeapAllocations();
27     struct avl_tree* tree = AvlTreeCreate();
28     assert(AvlTreeSize(tree) == 0);
29     AvlTreeInsert(tree, (void*) 5);
30     assert(AvlTreeSize(tree) == 1);
31     assert(AvlTreeContains(tree, (void*) 5));
32     assert(!AvlTreeContains(tree, (void*) 7));
33     assert(AvlTreeGet(tree, (void*) 5) == (void*) 5);
34     AvlTreeInsert(tree, (void*) 7);
35     AvlTreeInsert(tree, (void*) 9);
36     AvlTreeInsert(tree, (void*) 11);
37     AvlTreeInsert(tree, (void*) 8);
38     AvlTreeInsert(tree, (void*) 6);
39     assert(AvlTreeSize(tree) == 6);
40     AvlTreeInsert(tree, (void*) 4);
41     assert(AvlTreeSize(tree) == 7);
42     AvlTreeDelete(tree, (void*) 5);
43     assert(AvlTreeSize(tree) == 6);
44     assert(AvlTreeContains(tree, (void*) 4));
45     assert(!AvlTreeContains(tree, (void*) 5));
46     assert(AvlTreeContains(tree, (void*) 6));
47     assert(AvlTreeContains(tree, (void*) 7));
48     AvlTreeDestroy(tree);
49
50     /*
51     * Ensure there aren't any memory leaks.
52     */
53     assert(DbgGetOutstandingHeapAllocations() == heap_allocations);
54 }
55
56
57 void RegisterTfwAVLTreeTests(void) {
58     RegisterTfwTest("AVL trees (basic tests)", TFW_SP_AFTER_HEAP, AVLTreeBasic, PANIC_UNIT_TEST_OK, 0);
59 }
60
61 #endif
62

```

File: ./debug/tests/priorityqueue.c

```

1
2 #include <debug.h>
3 #include <assert.h>
4 #include <panic.h>
5 #include <string.h>
6 #include <log.h>
7 #include <arch.h>
8 #include <priorityqueue.h>
9 #include <stdlib.h>
10
11 #ifndef NDEBBUG
12
13 static void PQInsertionAndDeletionTest(bool max) {
14     struct priority_queue* queue = PriorityQueueCreate(100, max, 4);
15     assert(PriorityQueueGetCapacity(queue) == 100);
16     int elem;
17     for (int i = 0; i < 100; ++i) {
18         elem = i;
19         PriorityQueueInsert(queue, &elem, i * 3);
20         assert(PriorityQueueGetUsedSize(queue) == i + 1);
21     }
22
23     for (int i = 0; i < 100; ++i) {
24         struct priority_queue_result res = PriorityQueuePeek(queue);
25         int* d = res.data;
26         if (max) {
27             assert(*d == 99 - i);
28             assert((int) res.priority == (99 - i) * 3);
29         } else {
30             assert(*d == i);
31             assert((int) res.priority == i * 3);
32         }
33     }
34 }
35
36
37 void RegisterTfwPriorityQueueTests(void) {
38     RegisterTfwTest("Priority queue (basic tests)", TFW_SP_AFTER_HEAP, PQInsertionAndDeletionTest, PANIC_UNIT_TEST_OK, 0);
39 }
40
41 #endif
42

```

```

32     }
33     assert((int) PriorityQueueGetUsedSize(queue) == 100 - i);
34     PriorityQueuePop(queue);
35     assert((int) PriorityQueueGetUsedSize(queue) == 99 - i);
36 }
37 PriorityQueueDestroy(queue);
38 }
39
40 TFW_CREATE_TEST(PriorityQueueCombined) { TFW_IGNORE_UNUSED
41     PQInsertionAndDeletionTest(true);
42     PQInsertionAndDeletionTest(false);
43 }
44
45 TFW_CREATE_TEST(PriorityQueueStress) { TFW_IGNORE_UNUSED
46     int expected_size = 0;
47
48     srand(1);
49
50     struct priority_queue* queue = PriorityQueueCreate(100, true, 8);
51     assert(PriorityQueueGetCapacity(queue) == 100);
52
53     uint32_t data[2];
54
55     for (int i = 0; i < 1500000; ++i) {
56         int rng = rand();
57         for (int j = 0; j < 2; ++j) {
58             data[j] = rand();
59         }
60
61         if (rng % 3 == expected_size < 100) {
62             PriorityQueueInsert(queue, data, rng % 10000);
63             ++expected_size;
64
65         } else if (rng % 3 == 0 == expected_size > 1) {
66             struct priority_queue_result r1 = PriorityQueuePeek(queue);
67             PriorityQueuePop(queue);
68             struct priority_queue_result r2 = PriorityQueuePeek(queue);
69             PriorityQueuePop(queue);
70             assert(r1.priority == r2.priority);
71             expected_size -= 2;
72
73         } else {
74             --i;
75         }
76
77         assert(PriorityQueueGetUsedSize(queue) == expected_size);
78     }
79
80     uint64_t prev = 99999999;
81     while (PriorityQueueGetUsedSize(queue) > 0) {
82         struct priority_queue_result r1 = PriorityQueuePeek(queue);
83         PriorityQueuePop(queue);
84         assert(r1.priority <= prev);
85         prev = r1.priority;
86     }
87
88     PriorityQueueDestroy(queue);
89 }
90
91 TFW_CREATE_TEST(PriorityQueueInsertWhenFull) { TFW_IGNORE_UNUSED
92     int i = 0;
93     struct priority_queue* queue = PriorityQueueCreate(1, true, 4);
94     PriorityQueueInsert(queue, &i, 0);
95     PriorityQueueInsert(queue, &i, 0);
96 }
97
98
99 TFW_CREATE_TEST(PriorityQueuePeekWhenEmpty) { TFW_IGNORE_UNUSED
100     PriorityQueuePeek(PriorityQueueCreate(1, true, 4));
101 }
102
103 TFW_CREATE_TEST(PriorityQueuePopWhenEmpty) { TFW_IGNORE_UNUSED
104     PriorityQueuePop(PriorityQueueCreate(1, true, 4));
105 }
106
107 TFW_CREATE_TEST(PriorityQueueStangeSizes1) { TFW_IGNORE_UNUSED
108     PriorityQueueCreate(0, true, 4);
109 }
110
111 TFW_CREATE_TEST(PriorityQueueStangeSizes2) { TFW_IGNORE_UNUSED
112     PriorityQueueCreate(-1, true, 4);
113 }
114
115 TFW_CREATE_TEST(PriorityQueueStangeSizes3) { TFW_IGNORE_UNUSED
116     PriorityQueueCreate(1, true, 0);
117 }
118
119 TFW_CREATE_TEST(PriorityQueueStangeSizes4) { TFW_IGNORE_UNUSED
120     PriorityQueueCreate(1, true, -5);
121 }
122
123 void RegisterTfwPriorityQueueTests(void) {
124     RegisterTfwTest("Priority queues (general tests)", TFW_SP_AFTER_HEAP, PriorityQueueCombined, PANIC_UNIT_TEST_OK, 0);
125     RegisterTfwTest("Priority queues (stress tests)", TFW_SP_AFTER_HEAP, PriorityQueueStress, PANIC_UNIT_TEST_OK, 0);
126     RegisterTfwTest("Priority queues (insert when full)", TFW_SP_AFTER_HEAP, PriorityQueueInsertWhenFull, PANIC_PRIORITY_QUEUE, 0);
127     RegisterTfwTest("Priority queues (peek when empty)", TFW_SP_AFTER_HEAP, PriorityQueuePeekWhenEmpty, PANIC_PRIORITY_QUEUE, 0);
128     RegisterTfwTest("Priority queues (pop when empty)", TFW_SP_AFTER_HEAP, PriorityQueuePopWhenEmpty, PANIC_PRIORITY_QUEUE, 0);
129     RegisterTfwTest("Priority queues (strange sizes, 1)", TFW_SP_AFTER_HEAP, PriorityQueueStangeSizes1, PANIC_ASSERTION_FAILURE, 0);
130     RegisterTfwTest("Priority queues (strange sizes, 2)", TFW_SP_AFTER_HEAP, PriorityQueueStangeSizes2, PANIC_ASSERTION_FAILURE, 0);
131     RegisterTfwTest("Priority queues (strange sizes, 3)", TFW_SP_AFTER_HEAP, PriorityQueueStangeSizes3, PANIC_ASSERTION_FAILURE, 0);
132     RegisterTfwTest("Priority queues (strange sizes, 4)", TFW_SP_AFTER_HEAP, PriorityQueueStangeSizes4, PANIC_ASSERTION_FAILURE, 0);
133 }
134
135 #endif
136

```

File: ./debug/tests/init.c

```
1
2 #include <debug.h>
3 #include <assert.h>
4 #include <panic.h>
5 #include <string.h>
6 #include <log.h>
7 #include <arch.h>
8 #include <physical.h>
9
10 #ifndef NDEBUG
11
12 TFW_CREATE_TEST(BootSuccessful) { TFW_IGNORE_UNUSED
13
14 }
15
16 void RegisterTfwInitTests(void) {
17     RegisterTfwTest("Is boot successful", TFW_SP_ALL_CLEAR, BootSuccessful, PANIC_UNIT_TEST_OK, 0);
18 }
19
20 #endif
21 void RegisterTfwAVLTreeTests(void);
22 void RegisterTfwPriorityQueueTests(void);
```

File: ./debug/tests/physical.c

```

1
2 #include <debug.h>
3 #include <assert.h>
4 #include <panic.h>
5 #include <string.h>
6 #include <log.h>
7 #include <arch.h>
8 #include <physical.h>
9 #include <stdlib.h>
10
11 #ifndef NDEBUG
12
13 TFW_CREATE_TEST(IsPageAligned) ( TFW_IGNORE_UNUSED
14     assert(AllocPhys() % ARCH_PAGE_SIZE == 0);
15 )
16
17 TFW_CREATE_TEST(DeallocationChecksForPageAlignment) ( TFW_IGNORE_UNUSED
18     size_t p = AllocPhys();
19     DeallocPhys(p + context);
20 )
21
22 TFW_CREATE_TEST(DoubleDeallocationFails) ( TFW_IGNORE_UNUSED
23     size_t p = AllocPhys();
24     DeallocPhys(p);
25     DeallocPhys(p);
26 )
27
28 TFW_CREATE_TEST(SanityCheck) ( TFW_IGNORE_UNUSED
29     AllocPhys();
30     AllocPhys();
31 )
32
33 TFW_CREATE_TEST(BasicAllocationTest) ( TFW_IGNORE_UNUSED
34     size_t a = AllocPhys();
35     size_t b = AllocPhys();
36     size_t c = AllocPhys();
37     assert(a != b);
38     assert(a != c);
39     assert(b != c);
40 )
41
42 TFW_CREATE_TEST(BasicDeallocationTest) ( TFW_IGNORE_UNUSED
43     size_t a = AllocPhys();
44     size_t b = AllocPhys();
45     DeallocPhys(a);
46     DeallocPhys(b);
47     DeallocPhys(AllocPhys());
48 )
49
50 TFW_CREATE_TEST(StressTest) ( TFW_IGNORE_UNUSED
51     size_t frames[512];
52     inline_memset(frames, 0, sizeof(frames));
53     int allocated = 0;
54
55     srand(context * 1234 + 12);
56
57     // context 0: 40,000 (should be around 200ms)
58     // context 1: 250,000 (should be around 1s)
59     // context 2: 3,640,000 (NIGHTLY) (should be around 18s)
60     // context 3: 24,070,000 (NIGHTLY) (should be around 120s)
61
62     int limit = 10000 * (context * context * context * 3 + 2) * (context * context * 3 + 2);
63     for (int i = 0; i < limit; ++i) {
64         /*
65          * This all gets a bit dodgy if it's too much higher than 400 on a 4MB system, as it will
66          * eventually need to evict pages, but because we haven't actually mapped any of them into
67          * virtual memory, we just hang.
68          */
69         if (allocated < (300 + rand() % 80)) {
70             size_t f = AllocPhys();
71             for (int j = 0; j < 512; ++j) {
72                 if (frames[j] == 0) {
73                     Panic(PANIC_MANUALLY_INITIATED);
74                 }
75             }
76             for (int j = 0; j < 512; ++j) {
77                 if (frames[j] == 0) {
78                     ++allocated;
79                     frames[j] = f;
80                     break;
81                 }
82             }
83         } else {
84             while (allocated > 0) {
85                 for (int j = 0; j < 512; ++j) {
86                     int k = rand() % 256;
87                     if (frames[k] != 0) {
88                         DeallocPhys(frames[k]);
89                         --allocated;
90                         frames[k] = 0;
91                         break;
92                     }
93                 }
94                 if (rand() % 17 == 0) {
95                     break;
96                 }
97             }
98         }
99     }
100 }
101
102 TFW_CREATE_TEST(ContiguousAllocationRequiresStackAllocator) ( TFW_IGNORE_UNUSED
103     assert(AllocPhysContiguous(ARCH_PAGE_SIZE, 0, 0, 0) == 0);
104 )
105
106 void RegisterTfwPhysTests(void) {
107     RegisterTfwTest("Is AllocPhys sane", TFW_SP_AFTER_PHYS, SanityCheck, PANIC_UNIT_TEST_OK, 0);
108     RegisterTfwTest("Basic AllocPhys test (bitmap)", TFW_SP_AFTER_PHYS, BasicAllocationTest, PANIC_UNIT_TEST_OK, 0);
109     RegisterTfwTest("Basic AllocPhys test (stack)", TFW_SP_AFTER_PHYS_REINIT, BasicAllocationTest, PANIC_UNIT_TEST_OK, 0);
110     RegisterTfwTest("Basic DeallocPhys test (bitmap)", TFW_SP_AFTER_PHYS, BasicDeallocationTest, PANIC_UNIT_TEST_OK, 0);
111     RegisterTfwTest("Basic DeallocPhys test (stack)", TFW_SP_AFTER_PHYS_REINIT, BasicDeallocationTest, PANIC_UNIT_TEST_OK, 0);
112     RegisterTfwTest("AllocPhys and DeallocPhys stress test (bitmap 1)", TFW_SP_AFTER_PHYS, StressTest, PANIC_UNIT_TEST_OK, 0);
113     RegisterTfwTest("AllocPhys and DeallocPhys stress test (stack 1)", TFW_SP_AFTER_PHYS_REINIT, StressTest, PANIC_UNIT_TEST_OK, 0);
114     RegisterTfwTest("AllocPhys and DeallocPhys stress test (stack 2)", TFW_SP_AFTER_PHYS_REINIT, StressTest, PANIC_UNIT_TEST_OK, 1);
115     RegisterTfwTest("AllocPhys and DeallocPhys stress test (stack 3)", TFW_SP_AFTER_PHYS_REINIT, StressTest, PANIC_UNIT_TEST_OK, 2);
116     RegisterTfwTest("AllocPhys and DeallocPhys stress test (stack 4)", TFW_SP_AFTER_PHYS_REINIT, StressTest, PANIC_UNIT_TEST_OK, 3);
117     RegisterTfwTest("AllocPhys returns page aligned addresses", TFW_SP_AFTER_HEAP, IsPageAligned, PANIC_UNIT_TEST_OK, 0);
118     RegisterTfwTest("DeallocPhys only accepts page aligned addresses (1)", TFW_SP_AFTER_HEAP, DeallocationChecksForPageAlignment, PANIC_ASSERTION_FAILURE, 1);
119     RegisterTfwTest("DeallocPhys only accepts page aligned addresses (2)", TFW_SP_AFTER_HEAP, DeallocationChecksForPageAlignment, PANIC_ASSERTION_FAILURE, ARCH_P);
120     RegisterTfwTest("DeallocPhys checks for double allocation", TFW_SP_AFTER_HEAP, DoubleDeallocationFails, PANIC_ASSERTION_FAILURE, 0);
121
122     RegisterTfwTest("AllocPhysContiguous requires the stack allocator", TFW_SP_AFTER_PHYS, ContiguousAllocationRequiresStackAllocator, PANIC_UNIT_TEST_OK, 0);
123
124     // TODO: contiguous tests...
125 }
126
127 #endif
128

```

File: ./debug/tests/semaphore.c

```
1
2 #include <semaphore.h>
3 #include <debug.h>
4 #include <assert.h>
5 #include <panic.h>
6 #include <string.h>
7 #include <log.h>
8 #include <arch.h>
9 #include <irq.h>
10 #include <thread.h>
11 #include <errno.h>
12 #include <virtual.h>
13 #include <stdlib.h>
14
15 #ifndef NDEBBUG
16
17 static bool Thread1Ok = false;
18 static void Thread1(void* sem_) {
19     struct semaphore* sem = (struct semaphore*) sem_;
20
21     int res = AcquireSemaphore(sem, 1500);
22     Thread1Ok = true;
23     assert(res == ETIMEDOUT);
24
25     while (true) {
26         Schedule();
27     }
28 }
29
30 TFW_CREATE_TEST(SemaphoreTimeout1) { TFW_IGNORE_UNUSED
31     EXACT_IRQL(IRQL_STANDARD);
32
33     struct semaphore* sem = CreateSemaphore("", 1, 0);
34     AcquireSemaphore(sem, -1);
35
36     CreateThread(Thread1, (void*) sem, GetVas(), "");
37
38     SleepMilli(1000);
39     assert(!Thread1Ok);
40     SleepMilli(1000);
41     assert(Thread1Ok);
42 }
43
44 static void Thread2(void* sem_) {
45     struct semaphore* sem = (struct semaphore*) sem_;
46
47     int res = AcquireSemaphore(sem, 15000);
48     Thread1Ok = true;
49     assert(res == 0);
50
51     while (true) {
52         Schedule();
53     }
54 }
55
56 TFW_CREATE_TEST(SemaphoreTimeout2) { TFW_IGNORE_UNUSED
57     EXACT_IRQL(IRQL_STANDARD);
58
59     struct semaphore* sem = CreateSemaphore("", 1, 0);
60     AcquireSemaphore(sem, -1);
61
62     CreateThread(Thread2, (void*) sem, GetVas(), "");
63
64     SleepMilli(1000);
65     ReleaseSemaphore(sem);
66     SleepMilli(500);
67     assert(Thread1Ok);
68 }
69
70 static void Thread3(void* ignored) {
71     (void) ignored;
72
73     while (true) {
74         if (rand() % 10 == 0) {
75             for (int j = 0; j < rand() % 1000000; ++j) {
76                 DbgScreenPrintf("z");
77             }
78         }
79         SleepMilli(rand() % 150);
80     }
81 }
82
83 static void Thread3B(void* ignored) {
84     (void) ignored;
85
86     while (true) {
87         SleepMilli(rand() % 100);
88     }
89 }
90
91 static void Thread4(void* ignored) {
92     (void) ignored;
93
94     while (true) {
95         DbgScreenPrintf("y");
96     }
97 }
98
99 static struct semaphore* sems[20];
100
101 void Thread5(void* delay) {
102     int loops = 0;
103
104     while (true) {
105         int a = rand() % 20;
106         int b = rand() % 8;
107         int c = rand() % 3;
108         if (a == b || a == c || b == c) {
109             continue;
110         }
111
112         int ares = AcquireSemaphore(sems[a], (rand() % 10) * (rand() % 10));
113         int bres = AcquireSemaphore(sems[b], (rand() % 10) * (rand() % 10) * 2);
114         int cres = AcquireSemaphore(sems[c], (rand() % 10) * (rand() % 10) * 3);
115         if (delay != NULL) {
116             SleepMilli(5);
117         }
118         if (ares == 0) {
119             ReleaseSemaphore(sems[a]);
120         }
121         if (cres == 0) {
122             ReleaseSemaphore(sems[c]);
123         }
124         if (bres == 0) {
```



```

125         ReleaseSemaphore(sems.b);
126     }
127     ++loops;
128     if (rand() % 3 == 0) {
129         DbgScreenPrintf("%d,", loops);
130     }
131 }
132
133
134 TFW_CREATE_TEST(SchedulerHeartAttack) { TFW_IGNORE_UNUSED
135     EXACT__IRQL(IRQL_STANDARD);
136
137     for (int i = 0; i < 20; ++i) {
138         sems[i] = CreateSemaphore("", rand() % 3 + 1, 0);
139     }
140
141     for (int i = 0; i < 10; ++i) {
142         CreateThread(Thread3, NULL, GetVas(), "");
143     }
144     for (int i = 0; i < 10; ++i) {
145         CreateThread(Thread3B, NULL, GetVas(), "");
146     }
147     for (int i = 0; i < 5; ++i) {
148         CreateThread(Thread4, NULL, GetVas(), "");
149     }
150     for (int i = 0; i < 100; ++i) {
151         CreateThread(Thread5, context % 2 == 0 ? NULL : ((void*) 1), GetVas(), "");
152     }
153
154     for (int i = 0; i < (int) context; ++i) {
155         SleepMilli(1000);
156     }
157 }
158
159 void RegisterTfwSemaphoreTests(void) {
160     RegisterTfwTest("Semaphores with timeouts can be woken via timeout", TFW_SP_ALL_CLEAR, SemaphoreTimeout1, PANIC_UNIT_TEST_OK, 0);
161     RegisterTfwTest("Semaphores with timeouts can be woken via release", TFW_SP_ALL_CLEAR, SemaphoreTimeout2, PANIC_UNIT_TEST_OK, 0);
162     RegisterTfwTest("Scheduler stress test with semaphores (1)", TFW_SP_ALL_CLEAR, SchedulerHeartAttack, PANIC_UNIT_TEST_OK, 3);
163     RegisterTfwTest("Scheduler stress test with semaphores (2)", TFW_SP_ALL_CLEAR, SchedulerHeartAttack, PANIC_UNIT_TEST_OK, 4);
164     RegisterNightlyTfwTest("Scheduler stress test with semaphores (3)", TFW_SP_ALL_CLEAR, SchedulerHeartAttack, PANIC_UNIT_TEST_OK, 80);
165     RegisterNightlyTfwTest("Scheduler stress test with semaphores (4)", TFW_SP_ALL_CLEAR, SchedulerHeartAttack, PANIC_UNIT_TEST_OK, 125);
166     RegisterNightlyTfwTest("Scheduler stress test with semaphores (5)", TFW_SP_ALL_CLEAR, SchedulerHeartAttack, PANIC_UNIT_TEST_OK, 750);
167     RegisterNightlyTfwTest("Scheduler stress test with semaphores (6)", TFW_SP_ALL_CLEAR, SchedulerHeartAttack, PANIC_UNIT_TEST_OK, 755);
168 }
169
170 #endif

```

File: ./debug/tests/irq.c

```

1
2 #include <debug.h>
3 #include <assert.h>
4 #include <panic.h>
5 #include <string.h>
6 #include <log.h>
7 #include <arch.h>
8 #include <irq.h>
9
10 #ifndef NDEBBUG
11
12 static int counter = 0;
13
14 static void defer_me(void* context) {
15     counter = (size_t) context;
16 }
17
18 static void defer_me_2(void* context) {
19     counter += (size_t) context;
20 }
21
22 static void defer_me_3(void* context) {
23     int irq1 = RaiseIrql(IRQL_HIGH);
24     LowerIrql(irq1);
25
26     counter += (size_t) context;
27 }
28
29 static void internal_def4_1(void* context) {
30     (void) context;
31
32     assert(counter == 20);
33     counter = 30;
34 }
35
36 static void internal_def4_2(void* context) {
37     (void) context;
38
39     assert(counter == 40);
40     counter = 50;
41 }
42
43 static void defer_me_4(void* context) {
44     (void) context;
45
46     assert(counter == 0);
47
48     int irq1 = RaiseIrql(IRQL_HIGH);
49     DeferUntilIrql(IRQL_TIMER, internal_def4_1, NULL);
50     DeferUntilIrql(IRQL_PAGE_FAULT, internal_def4_2, NULL);
51     counter = 20;
52     LowerIrql(irq1);
53     assert(counter == 30);
54 }
55
56 static void defer_me_5(void* context) {
57     (void) context;
58
59     assert(counter == 30);
60     counter = 40;
61 }
62
63 static void defer_me_6(void* context) {
64     (void) context;
65     int v = (size_t) context;
66     assert(counter == v - 1);
67     counter = v;
68 }
69
70 TFW_CREATE_TEST(RaiseLowerTest) { TFW_IGNORE_UNUSED
71     int irq1_a = GetIrql();
72     assert(irq1_a == IRQL_STANDARD);
73     int irq1_b = RaiseIrql(IRQL_SCHEDULER);
74     assert(irq1_a == irq1_b);
75     assert(GetIrql() == IRQL_SCHEDULER);
76     int irq1_c = RaiseIrql(IRQL_HIGH);
77     assert(irq1_c == IRQL_SCHEDULER);
78     assert(GetIrql() == IRQL_HIGH);

```

```

79     LowerIrql(irql_c);
80     assert(GetIrql() == IRQL_SCHEDULER);
81     LowerIrql(irql_a);
82     assert(GetIrql() == irql_a);
83 }
84
85 TFW_CREATE_TEST(DeferRunsImmediatelyAtLevel) { TFW_IGNORE_UNUSED
86     EXACT_IRQL(IRQL_STANDARD);
87     counter = 0;
88     DeferUntilIrql(IRQL_STANDARD, defer_me, (void*) 1);
89     assert(counter == 1);
90     RaiseIrql(IRQL_SCHEDULER);
91     DeferUntilIrql(IRQL_SCHEDULER, defer_me, (void*) 2);
92     assert(counter == 2);
93     RaiseIrql(IRQL_HIGH);
94     DeferUntilIrql(IRQL_HIGH, defer_me, (void*) 3);
95     assert(counter == 3);
96 }
97
98 TFW_CREATE_TEST(DeferDoesntWorkBeforeHeap) { TFW_IGNORE_UNUSED
99     EXACT_IRQL(IRQL_STANDARD);
100     counter = 0;
101     DeferUntilIrql(IRQL_STANDARD, defer_me, (void*) 1);
102     assert(counter == 1);
103     RaiseIrql(IRQL_SCHEDULER);
104     DeferUntilIrql(IRQL_STANDARD, defer_me, (void*) 2);
105     assert(counter == 1);
106     LowerIrql(IRQL_SCHEDULER);
107     assert(counter == 1);
108 }
109
110 TFW_CREATE_TEST(DeferWorksNormally) { TFW_IGNORE_UNUSED
111     EXACT_IRQL(IRQL_STANDARD);
112     counter = 0;
113     DeferUntilIrql(IRQL_STANDARD, defer_me, (void*) 1);
114     assert(counter == 1);
115     RaiseIrql(IRQL_SCHEDULER);
116     DeferUntilIrql(IRQL_STANDARD, defer_me, (void*) 2);
117     assert(counter == 1);
118     LowerIrql(IRQL_STANDARD);
119     assert(counter == 2);
120 }
121
122 TFW_CREATE_TEST(DeferWorksThroughLevelsInOrder) { TFW_IGNORE_UNUSED
123     EXACT_IRQL(IRQL_STANDARD);
124     counter = 0;
125
126     RaiseIrql(IRQL_HIGH);
127     DeferUntilIrql(IRQL_HIGH, defer_me_6, (void*) 1);
128     assert(counter == 1);
129     DeferUntilIrql(IRQL_TIMER, defer_me_6, (void*) 2);
130     assert(counter == 1);
131     DeferUntilIrql(IRQL_DRIVER, defer_me_6, (void*) 3);
132     assert(counter == 1);
133     DeferUntilIrql(IRQL_SCHEDULER, defer_me_6, (void*) 4);
134     assert(counter == 1);
135     DeferUntilIrql(IRQL_PAGE_FAULT, defer_me_6, (void*) 5);
136     assert(counter == 1);
137     DeferUntilIrql(IRQL_STANDARD, defer_me_6, (void*) 6);
138     assert(counter == 1);
139
140     LowerIrql(IRQL_STANDARD);
141     assert(counter == 6);
142 }
143
144
145 TFW_CREATE_TEST(DeferWorksThroughLevels) { TFW_IGNORE_UNUSED
146     EXACT_IRQL(IRQL_STANDARD);
147     counter = 0;
148
149     RaiseIrql(IRQL_HIGH);
150     DeferUntilIrql(IRQL_HIGH, defer_me_2, (void*) 1);
151     assert(counter == 1);
152     DeferUntilIrql(IRQL_TIMER, defer_me_2, (void*) 2);
153     assert(counter == 1);
154     DeferUntilIrql(IRQL_DRIVER, defer_me_2, (void*) 4);
155     assert(counter == 1);
156     DeferUntilIrql(IRQL_SCHEDULER, defer_me_2, (void*) 8);
157     assert(counter == 1);
158     DeferUntilIrql(IRQL_PAGE_FAULT, defer_me_2, (void*) 16);
159     assert(counter == 1);
160     DeferUntilIrql(IRQL_STANDARD, defer_me_2, (void*) 32);
161     assert(counter == 1);
162
163     LowerIrql(IRQL_STANDARD);
164     assert(counter == 63);
165 }
166
167 TFW_CREATE_TEST(DeferMultipleAtSameLevel) { TFW_IGNORE_UNUSED
168     EXACT_IRQL(IRQL_STANDARD);
169     counter = 0;
170
171     RaiseIrql(IRQL_HIGH);
172     DeferUntilIrql(IRQL_TIMER, defer_me_2, (void*) 1);
173     assert(counter == 0);
174     DeferUntilIrql(IRQL_TIMER, defer_me_2, (void*) 2);
175     assert(counter == 0);
176     DeferUntilIrql(IRQL_TIMER, defer_me_2, (void*) 4);
177     assert(counter == 0);
178
179     LowerIrql(IRQL_STANDARD);
180     assert(counter == 7);
181 }
182
183 TFW_CREATE_TEST(DeferDoesntRunLowHandlers) { TFW_IGNORE_UNUSED
184     EXACT_IRQL(IRQL_STANDARD);
185     counter = 0;
186
187     RaiseIrql(IRQL_HIGH);
188     DeferUntilIrql(IRQL_HIGH, defer_me_2, (void*) 1);
189     assert(counter == 1);
190     DeferUntilIrql(IRQL_TIMER, defer_me_2, (void*) 2);
191     assert(counter == 1);
192     DeferUntilIrql(IRQL_DRIVER, defer_me_2, (void*) 4);
193     assert(counter == 1);
194     DeferUntilIrql(IRQL_SCHEDULER, defer_me_2, (void*) 8);
195     assert(counter == 1);
196     DeferUntilIrql(IRQL_PAGE_FAULT, defer_me_2, (void*) 16);
197     assert(counter == 1);
198     DeferUntilIrql(IRQL_STANDARD, defer_me_2, (void*) 32);
199     assert(counter == 1);
200
201     LowerIrql(IRQL_SCHEDULER);
202     assert(counter == 15);
203 }
204
205 TFW_CREATE_TEST(DeferWorksThroughLevelsStepping) { TFW_IGNORE_UNUSED
206     EXACT_IRQL(IRQL_STANDARD);
207     counter = 0;
208

```

```

209 RaiseIrql(IRQL_HIGH);
210 DeferUntilIrql(IRQL_HIGH, defer_me_2, (void*) 1);
211 assert(counter == 1);
212 DeferUntilIrql(IRQL_TIMER, defer_me_2, (void*) 2);
213 assert(counter == 1);
214 DeferUntilIrql(IRQL_TIMER, defer_me_2, (void*) 4);
215 assert(counter == 1);
216 DeferUntilIrql(IRQL_DRIVER, defer_me_2, (void*) 8);
217 assert(counter == 1);
218 DeferUntilIrql(IRQL_SCHEDULER, defer_me_2, (void*) 16);
219 assert(counter == 1);
220 DeferUntilIrql(IRQL_PAGE_FAULT, defer_me_2, (void*) 32);
221 assert(counter == 1);
222 DeferUntilIrql(IRQL_STANDARD, defer_me_2, (void*) 64);
223 assert(counter == 1);
224
225 LowerIrql(IRQL_DRIVER);
226 assert(counter == 15);
227 RaiseIrql(IRQL_HIGH);
228 LowerIrql(IRQL_DRIVER);
229 assert(counter == 15);
230 LowerIrql(IRQL_PAGE_FAULT);
231 assert(counter == 65);
232 LowerIrql(IRQL_STANDARD);
233 assert(counter == 127);
234
235
236 TFW_CREATE_TEST(DeferWithLoweringInHandler) ( TFW_IGNORE_UNUSED
237 EXACT_IRQL(IRQL_STANDARD);
238 counter = 0;
239
240 RaiseIrql(IRQL_SCHEDULER);
241 DeferUntilIrql(IRQL_PAGE_FAULT, defer_me_3, (void*) 55);
242 assert(counter == 0);
243 LowerIrql(IRQL_STANDARD);
244 assert(counter == 55);
245
246
247 TFW_CREATE_TEST(DeferWithDeferringInHandler) ( TFW_IGNORE_UNUSED
248 EXACT_IRQL(IRQL_STANDARD);
249 counter = 0;
250
251 RaiseIrql(IRQL_HIGH);
252 DeferUntilIrql(IRQL_DRIVER, defer_me_4, NULL);
253 DeferUntilIrql(IRQL_SCHEDULER, defer_me_5, NULL);
254 assert(counter == 0);
255 LowerIrql(IRQL_STANDARD);
256 assert(counter == 50);
257
258
259 void RegisterTfwIrqlTests(void) {
260 RegisterTfwTest("RaiseIrql, LowerIrql and GetIrql work", TFW_SP_AFTER_HEAP, RaiseLowerTest, PANIC_UNIT_TEST_OK, 0);
261 RegisterTfwTest("DeferUntilIrql gets run immediately at level", TFW_SP_AFTER_HEAP, DeferRunsImmediatelyAtLevel, PANIC_UNIT_TEST_OK, 0);
262 RegisterTfwTest("DeferUntilIrql gets run on level lowering", TFW_SP_AFTER_HEAP, DeferWorksNormally, PANIC_UNIT_TEST_OK, 0);
263 RegisterTfwTest("DeferUntilIrql defers get ignored before heap", TFW_SP_AFTER_PHYS, DeferDoesntWorkBeforeHeap, PANIC_UNIT_TEST_OK, 0);
264 RegisterTfwTest("DeferUntilIrql runs multiple handlers at same level", TFW_SP_AFTER_HEAP, DeferMultipleAtSameLevel, PANIC_UNIT_TEST_OK, 0);
265 RegisterTfwTest("DeferUntilIrql runs multiple handlers at different levels (1)", TFW_SP_AFTER_HEAP, DeferWorksThroughLevels, PANIC_UNIT_TEST_OK, 0);
266 RegisterTfwTest("DeferUntilIrql runs multiple handlers at different levels (2)", TFW_SP_AFTER_HEAP, DeferWorksThroughLevelsStepping, PANIC_UNIT_TEST_OK, 0);
267 RegisterTfwTest("DeferUntilIrql runs multiple handlers in the correct order", TFW_SP_AFTER_HEAP, DeferWorksThroughLevelsInOrder, PANIC_UNIT_TEST_OK, 0);
268 RegisterTfwTest("DeferUntilIrql doesn't run handlers below current level", TFW_SP_AFTER_HEAP, DeferDoesntRunLowHandlers, PANIC_UNIT_TEST_OK, 0);
269 RegisterTfwTest("DeferUntilIrql when handler calls LowerIrql", TFW_SP_AFTER_HEAP, DeferWithLoweringInHandler, PANIC_UNIT_TEST_OK, 0);
270 RegisterTfwTest("DeferUntilIrql when handler calls DeferUntilIrql", TFW_SP_AFTER_HEAP, DeferWithDeferringInHandler, PANIC_UNIT_TEST_OK, 0);
271
272
273 #endif

```

File: ./debug/tests/wait.c

```

1
2 #include <semaphore.h>
3 #include <debug.h>
4 #include <assert.h>
5 #include <panic.h>
6 #include <string.h>
7 #include <log.h>
8 #include <arch.h>
9 #include <irq.h>
10 #include <thread.h>
11 #include <process.h>
12 #include <errno.h>
13 #include <virtual.h>
14 #include <stdlib.h>
15
16 #ifndef NDEBUG
17
18 static void SecondProcessThread(void* arg) {
19     int status = (int) (size_t) arg;
20     SleepMilli(1000);
21     KillProcess(status);
22 }
23
24 static void ZombieProcess(void*) {
25     KillProcess(99);
26 }
27
28 static bool ok = false;
29
30 static void InitialProcessThread1(void*) {
31     struct process* child1 = CreateProcessWithEntryPoint(1, SecondProcessThread, (void*) (size_t) 111);
32     SleepMilli(300);
33     struct process* child2 = CreateProcessWithEntryPoint(1, SecondProcessThread, (void*) (size_t) 222);
34     SleepMilli(300);
35     struct process* child3 = CreateProcessWithEntryPoint(1, SecondProcessThread, (void*) (size_t) 333);
36     int retv;
37     WaitProcess(GetPid(child3), &retv, 0);
38     assert(retv == 333);
39     WaitProcess(GetPid(child2), &retv, 0);
40     assert(retv == 222);
41     WaitProcess(GetPid(child1), &retv, 0);
42     assert(retv == 111);
43     ok = true;
44 }
45
46 static void InitialProcessThread2(void*) {
47     pid_t c1pid = GetPid(CreateProcessWithEntryPoint(1, SecondProcessThread, (void*) (size_t) 111));
48     SleepMilli(500);
49     pid_t c2pid = GetPid(CreateProcessWithEntryPoint(1, SecondProcessThread, (void*) (size_t) 222));
50     SleepMilli(500);
51     pid_t c3pid = GetPid(CreateProcessWithEntryPoint(1, SecondProcessThread, (void*) (size_t) 333));
52
53     int retv;
54     pid_t pid = WaitProcess(-1, &retv, 0);
55     assert(retv == 111);
56     assert(pid == c1pid);
57
58     pid = WaitProcess(c3pid, &retv, 0);

```

```

59     assert(retv == 333);
60     assert(pid == c3pid);
61
62     pid = WaitProcess(-1, &retv, 0);
63     assert(retv == 222);
64     assert(pid == c2pid);
65
66     ok = true;
67 }
68
69 static void InitialProcessThread3(void*) {
70     pid_t zombie = GetPid(CreateProcessWithEntryPoint(1, ZombieProcess, NULL));
71     SleepMilli(500);
72
73     int retv;
74     pid_t pid = WaitProcess(-1, &retv, 0);
75     assert(retv == 99);
76     assert(pid == zombie);
77     ok = true;
78 }
79
80 static void InitialProcessThread4(void*) {
81     pid_t zombie = GetPid(CreateProcessWithEntryPoint(1, ZombieProcess, NULL));
82     SleepMilli(500);
83
84     int retv;
85     pid_t pid = WaitProcess(zombie, &retv, 0);
86     assert(retv == 99);
87     assert(pid == zombie);
88     ok = true;
89 }
90
91 static void InitialProcessThread5(void* mode_) {
92     size_t mode = (size_t) mode_;
93
94     pid_t pids[30];
95     for (int i = 0; i < 30; ++i) {
96         pids[i] = GetPid(CreateProcessWithEntryPoint(1, ZombieProcess, NULL));
97     }
98
99     int retv;
100     if (mode == 0) {
101         for (int i = 0; i < 30; ++i) {
102             WaitProcess(-1, &retv, 0);
103             assert(retv == 99);
104         }
105     } else if (mode == 1) {
106         for (int i = 0; i < 30; ++i) {
107             pid_t pid = WaitProcess(pids[i], &retv, 0);
108             assert(retv == 99);
109             assert(pid == pids[i]);
110         }
111     } else {
112         for (int i = 0; i < 30; ++i) {
113             pid_t pid = WaitProcess(pids[29 - i], &retv, 0);
114             assert(retv == 99);
115             assert(pid == pids[29 - i]);
116         }
117     }
118
119     ok = true;
120 }
121
122 TFW_CREATE_TEST(BasicWaitTest) ( TFW_IGNORE_UNUSED
123     EXACT_IRQL(IRQL_STANDARD);
124     CreateProcessWithEntryPoint(0, InitialProcessThread1, NULL);
125     SleepMilli(2000);
126     assert(ok);
127 )
128
129 TFW_CREATE_TEST(WaitTestWithNeg1) ( TFW_IGNORE_UNUSED
130     EXACT_IRQL(IRQL_STANDARD);
131     CreateProcessWithEntryPoint(0, InitialProcessThread2, NULL);
132     SleepMilli(3000);
133     assert(ok);
134 )
135
136 TFW_CREATE_TEST(WaitOnZombieTest1) ( TFW_IGNORE_UNUSED
137     EXACT_IRQL(IRQL_STANDARD);
138     CreateProcessWithEntryPoint(0, InitialProcessThread3, NULL);
139     SleepMilli(1000);
140     assert(ok);
141 )
142
143 TFW_CREATE_TEST(WaitOnZombieTest2) ( TFW_IGNORE_UNUSED
144     EXACT_IRQL(IRQL_STANDARD);
145     CreateProcessWithEntryPoint(0, InitialProcessThread4, NULL);
146     SleepMilli(1000);
147     assert(ok);
148 )
149
150 TFW_CREATE_TEST(WaitOnManyTest1) ( TFW_IGNORE_UNUSED
151     EXACT_IRQL(IRQL_STANDARD);
152     CreateProcessWithEntryPoint(0, InitialProcessThread5, (void*) (size_t) 0);
153     SleepMilli(6000);
154     assert(ok);
155 )
156
157 TFW_CREATE_TEST(WaitOnManyTest2) ( TFW_IGNORE_UNUSED
158     EXACT_IRQL(IRQL_STANDARD);
159     CreateProcessWithEntryPoint(0, InitialProcessThread5, (void*) (size_t) 1);
160     SleepMilli(6000);
161     assert(ok);
162 )
163
164 TFW_CREATE_TEST(WaitOnManyTest3) ( TFW_IGNORE_UNUSED
165     EXACT_IRQL(IRQL_STANDARD);
166     CreateProcessWithEntryPoint(0, InitialProcessThread5, (void*) (size_t) 2);
167     SleepMilli(6000);
168     assert(ok);
169 )
170
171
172 void RegisterTfwWaitTests(void) {
173     RegisterTfwTest("WaitProcess works (explicit ids)", TFW_SP_ALL_CLEAR, BasicWaitTest, PANIC_UNIT_TEST_OK, 0);
174     RegisterTfwTest("WaitProcess works (-1)", TFW_SP_ALL_CLEAR, WaitTestWithNeg1, PANIC_UNIT_TEST_OK, 0);
175     RegisterTfwTest("WaitProcess allows waiting on zombie (general)", TFW_SP_ALL_CLEAR, WaitOnZombieTest1, PANIC_UNIT_TEST_OK, 0);
176     RegisterTfwTest("WaitProcess allows waiting on zombie (explicit)", TFW_SP_ALL_CLEAR, WaitOnZombieTest2, PANIC_UNIT_TEST_OK, 0);
177     RegisterTfwTest("WaitProcess works when waiting on many (general)", TFW_SP_ALL_CLEAR, WaitOnManyTest1, PANIC_UNIT_TEST_OK, 0);
178     RegisterTfwTest("WaitProcess works when waiting on many (explicit, in order)", TFW_SP_ALL_CLEAR, WaitOnManyTest2, PANIC_UNIT_TEST_OK, 0);
179     RegisterTfwTest("WaitProcess works when waiting on many (explicit, reversed)", TFW_SP_ALL_CLEAR, WaitOnManyTest3, PANIC_UNIT_TEST_OK, 0);
180     // todo: stress tests
181 }
182
183 #endif

```