# ToaruOS at 5 Years

A Closer Look at a Hobbyist Operating System



#### Outline

Background & Design Goals

"What happens when..." - deep walkthrough

Demos

Future

Questions?

# Background

## Background

とあるOS - Pronounced "Toe Ah Roo Oh Es"

Means something along the lines of "A Specific OS" or "Some Such OS"

Parody of generic hobby OS names like "MyOS"

## Background

Project started in December, 2010

First code commits on 2011-01-15

First screenshot from 2011-01-20

ToAruOS 0.0.1]							
'lags : 0x0000004f Mem Lo	: 0x0000027f	Mem Hi:	0×0001fb80	Boot d:	0×8000ffff		
mdlin: 0x00108000 Mods	: 0x00000000	Addr :	0×00108000	Syms :	0×00000000		
yms : 0x00000000 Syms	: 0x00000000	Syms :	0×00000000	MMap :	0×00000090		
ddr : 0x00009000 Drives	: 0x00000000	Addr :	0×00000000	Config:	0x00000000		
.oader: 0x00000000 APM	: 0x00000000	<b>VBE Co:</b>	0×00000000	VBE Mo:	0×00000000		
JBE In: 0x00000000 UBE se	e: 0x00000000	VBE of:	0×00000000	VBE le:	0×00000000		
End multiboot raw data)							
tarted with:							
looted from: S							
39kB lower memory							
29920kB higher memory (126MB)							
esting colors							
= Roadmap =							
Paging							
Heap							
VFS							
Initial ramdisk							
Task switching							
User mode execution							
EXT2							
ernel is finished bootin	ıg. Press `q`	to produ	ice a page f	`ault.			

### Design Goals - Then

POSIX

x86

С

### **Overall Design**

Monolithic, modular

32-bit x86

Unix-like

Compositing window manager

### Publicity

github.com/showcases/open-source-operating-systems

phoronix.com/scan.php?page=news\_item&px=MTgzMTk

osnews.com/story/28038/ToAruOS\_hobby\_kernel\_and\_userspace

## April Fools: PonyOS



... you type "google.com" into your web browser and hit <enter>?

... you type "google.com" into your web browser and hit <enter>?

... you type `fetch http://toaruos.org/docs/talk.pdf` in a shell and hit <enter>?

### You are standing in an open field...

First, let's talk about the initial state of our system.

We'll assume we've already opened a terminal, the shell is running, and it is waiting for our input.

But what does that look like in code?

First off, let's talk about the shell.

ToaruOS's shell implements complex line editing using a function called rline, which you can see to the right.

Since we're waiting for input, rline is blocked in the call to fgetc on line 48.

```
28      int rline(char * buffer, int buf size, rline callbacks t * call
29
       /* Initialize context */
30
       rline context t context = {
31
           buffer,
           callbacks,
           Θ,
34
           buf size,
           Θ,
37
           0,
           Θ.
       };
41
       printf("\033[s");
42
       fflush(stdout);
43
44
       key event state t kbd state = \{0\};
45
       /* Read kevs */
47
       while ((context.collected < context.requested) && (!context</pre>
           uint32 t key sym = kbd key(&kbd state, fgetc(stdin));
48
49
           if (key sym == KEY NONE) continue;
           switch (key_sym) {
               case KEY CTRL C:
                    printf("^C\n");
                    context.buffer[0] = '\0';
54
                    return 0;
               case KEY CTRL R:
                    if (callbacks->rev search) {
                        callbacks->rev search(&context);
                        return context.collected;
59
                    continue;
               case KEY ARROW UP:
62
               case KEY CTRL P:
                    if (callbacks->key up) {
64
                        callbacks->kev un(&context):
```

Skipping over the details of the C library, this call to fgetc eventually turns into a system call, so we'll turn our attention to the system call handler and sys\_read in the kernel.

sys\_read takes three arguments: a file descriptor number, a pointer to a buffer, and a length of data to read. The kernel will use the current process's file descriptor table to find the right read method to run.

```
753 void syscall handler(struct regs * r) {
       if (r->eax >= num syscalls) {
755
            return:
756
       }
758
       uintptr t location = (uintptr t)syscalls[r->eax];
759
       if (!location) {
760
            return:
        }
762
763
       /* Update the syscall registers for this process */
764
       current_process->syscall_registers = r;
766
       /* Call the syscall function */
       scall func func = (scall func)location;
768
       uint32_t ret = func(r->ebx, r->ecx, r->edx, r->esi, r->edi);
769
770
       if ((current process->syscall registers == r) ||
                (location != (uintptr_t)&fork && location != (uintptr_t)&clone)) {
771
772
            r \rightarrow eax = ret;
773
774 }
776
58 static int sys_read(int fd, char * ptr, int len) {
59
        if (FD CHECK(fd)) {
60
            PTR VALIDATE(ptr);
            fs node t * node = FD ENTRY(fd);
62
            uint32_t out = read_fs(node, node->offset, len, (uint8_t *)ptr);
63
            node->offset += out;
            return (int)out;
65
        }
        return -1;
```

The file descriptor that our fgetc call is trying to read on is 0 standard input, which for the shell points to a pseudo terminal *device*, so our next stop is the *tty* driver, where we find read pty slave. We'll skip over some details of tty modes and note that read pty slave has called ring buffer read. TTYs in ToaruOS use a *ring buffer* to store data that has been written on one end but not read from the other.



So what's our ring buffer doing? We've asked to read from an empty ring buffer, so the ring buffer needs to wait until data is available to read. It calls sleep\_on to wait until then.

What will wake it up? Writing to the other end of the TTY, which will be done by the *terminal emulator*.

51	size t ring buffer read(ring buffer t * ring buffer, size t size, uint8 t * buffer) {		
52	size t collected = 0:		
53	while (collected == $0$ ) {		
54	<pre>spin_lock(ring buffer-&gt;lock);</pre>		
55	while (ring_buffer_unread(ring_buffer) > 0 && collected < size) {		
56	buffer[collected] = ring_buffer->buffer[ring_buffer->read_ptr];		
57	ring_buffer_increment_read(ring_buffer);		
58	collected++;		
59	}		
60	spin_unlock(ring_buffer->lock);		
61	wakeup_queue(ring_buffer->wait_queue_writers);		
62	if (collected == 0) {		
63	if (sleep_on(ring_buffer->wait_queue_readers) && ring_buffer->internal_stop) {		
64	ring_buffer->internal_stop = 0;		
65	break;		
66			
67			
68			
69	wakeup_queue(ring_buffer->wait_queue_writers);		
70	return collected;		
/1			
70			

As such, our next stop in the codebase is the terminal emulator - one of ToaruOS's earliest graphical applications.

The terminal is actually doing a lot of things, but we're particularly interested in how it will receive our key press when we first type 'f' and for that we need to look at the handle\_incoming method.

```
1121 void * handle incoming(void * garbage) {
1122
         while (!exit_application) {
1123
             yutani msg t * m = yutani poll(yctx);
1124
             if (m) {
1125
                 switch (m->type) {
1126
                      case YUTANI_MSG_KEY_EVENT:
1127
1128
                              struct yutani msg key event * ke
1129
                              int ret = (ke->event.action == KE
1130
                              key_event(ret, &ke->event);
1131
1132
                          break:
1133
                     case YUTANI MSG WINDOW FOCUS CHANGE:
1134
1135
                              struct yutani msg window focus ch
1136
                              yutani window t * win = hashmap q
1137
                              if (win) {
1138
                                  win->focused = wf->focused:
1139
                                  render decors();
1140
1141
1142
                          break:
1143
                     case YUTANI_MSG_SESSION_END:
1144
1145
                              kill(child pid, SIGKILL);
1146
                              exit application = 1;
1147
1148
                          break:
1149
                     case YUTANI_MSG_RESIZE_OFFER:
1150
1151
                              struct yutani msg window resize *
1152
                              resize finish(wr->width, wr->heig
1153
1154
                          break:
1155
                     case YUTANI MSG WINDOW MOUSE EVENT:
```

The terminal is waiting for events from the *window compositor*, which will send messages through a ToaruOS-specific messaging system called a *packet exchange*.

That means we need to look at another application to see how our key presses get to the shell, so next we'll look at the compositor.

```
1121 void * handle incoming(void * garbage) {
1122
         while (!exit_application) {
1123
             yutani_msg_t * m = yutani_poll(yctx);
1124
             if (m) {
1125
                 switch (m->type) {
1126
                      case YUTANI_MSG_KEY_EVENT:
1127
1128
                              struct yutani msg key event * ke
1129
                              int ret = (ke->event.action == KE
1130
                              key_event(ret, &ke->event);
1131
1132
                          break:
1133
                     case YUTANI_MSG_WINDOW_FOCUS_CHANGE:
1134
1135
                              struct yutani msg window focus ch
1136
                              yutani window t * win = hashmap q
1137
                              if (win) {
1138
                                  win->focused = wf->focused:
                                  render_decors();
1139
1140
1141
1142
                          break:
1143
                     case YUTANI_MSG_SESSION_END:
1144
1145
                              kill(child pid, SIGKILL);
1146
                              exit application = 1;
1147
1148
                          break:
1149
                     case YUTANI_MSG_RESIZE_OFFER:
1150
1151
                              struct yutani msg window resize *
1152
                              resize finish(wr->width, wr->heig
1153
1154
                          break:
1155
                     case YUTANI MSG WINDOW MOUSE EVENT:
```

The compositor performs many functions, but we'll focus on how it handles keyboard input for now.

The compositor runs a thread to read scancodes from the *keyboard device file* and turns them into messages which it sends back to another thread.

The read call on line 562 directly translates to a read system call that is very similar to the one from the shell.

```
546
     * Keyboard input thread
547
548
     * Reads the kernel keyboard device and converts
549
      key presses into event objects to send to the
550
     * core compositor.
552
553 void * keyboard input(void * garbage) {
        int kfd = open("/dev/kbd", 0 RDONLY);
554
555
        yutani t * y = yutani init();
556
        key event t event;
558
        key_event_state_t state = {0};
559
        while (1) {
            char buf[1];
            int r = read(kfd, buf, 1);
562
            if (r > 0) {
564
                kbd scancode(&state, buf[0], &event);
565
                yutani msg t * m = yutani msg build key event(0, &event, &state);
                int result = yutani_msg_send(y, m);
567
                free(m):
568
            }
569
570 }
```

#### Hurry up and wait...

So our shell, terminal, and compositor are all sitting around *waiting* for other things to happen.

Let's see how things unfold when we hit that first key: *f* 

As soon as we hit the key, a hardware interrupt triggers on the CPU. Long before our shell ever started running, the *PS/2* keyboard driver set up an interrupt handler, and so the CPU will call our keyboard\_handler function.

The keyboard driver is very simple and just fills a buffer with the incoming scancode, which will wake up the compositor.

```
37 static int keyboard_handler(struct regs *r) {
38     unsigned char scancode;
39     keyboard_wait();
40     scancode = inportb(KEY_DEVICE);
41     irq_ack(KEY_IRQ);
42
43     write_fs(keyboard_pipe, 0, 1, (uint8_t []){scancode});
44     return 1;
45 }
46
```

Back in the compositor, our read call has finished and we have a scancode. We turn that scancode into a key event message and send it to the compositor's main thread as a packet exchange message with yutani\_msg\_send.

```
546 /**
547
    * Keyboard input thread
548
549
     * Reads the kernel keyboard device and converts
550
     * core compositor.
553 void * keyboard input(void * garbage) {
554
        int kfd = open("/dev/kbd", 0 RDONLY);
555
        yutani_t * y = yutani_init();
        key event t event;
558
        key_event_state_t state = {0};
559
        while (1) {
            char buf[1];
            int r = read(kfd, buf, 1);
            if (r > 0) {
564
                kbd_scancode(&state, buf[0], &event);
565
                yutani msg t * m = yutani msg build key event(0, &event, &state);
                int result = yutani_msg_send(y, m);
                free(m):
            }
569
570 }
```

The compositor's main event handler thread is waiting to receive messages, both from the keyboard thread and from all of the client applications that are communicating with it.

It will see the key event and then pass it off to handle\_key\_event.

There, it will be checked against various window management keybindings and eventually passed off to the focused window's application.

```
2014
2015
             while (1) {
                    pex packet t * p = calloc(PACKET SIZE, 1);
2016
2017
                    pex listen(server, p);
2018
2019
                    yutani_msg_t * m = (yutani_msg_t *)p->data;
2020
2086
                   case YUTANI_MSG_KEY_EVENT:
2087
2088
                            /* XXX Verify this is from a valid device client */
                            struct yutani_msg_key_event * ke = (void *)m->data;
2089
2090
                            handle key event(yg, ke);
2091
2092
                       break:
.428 static void <mark>handle_key_event</mark>(yutani_globals_t * yg, struct yutani_msg_key_event * ke) {
429
        yutani_server_window_t * focused = get_focused(yg);
1430
       memcpy(&yg->kbd state, &ke->state, sizeof(key event state t));
        if (focused) {
432
           if ((ke->event.action == KEY_ACTION_DOWN) &&
1433
                (ke->event.modifiers & KEY MOD LEFT CTRL) &&
                (ke->event.modifiers & KEY_MOD_LEFT_SHIFT) &&
435
                (ke->event.keycode == 'z')) {
1436
               mark window(yg,focused);
               focused->rotation -= 5;
               mark window(yg,focused);
                return;
       /* Finally, send the key to the focused client. */
1584
       if (focused) {
          yutani msg t * response = yutani msg build key event(focused->wid, &ke->event, &ke->state)
          pex send(yg->server, focused->owner, response->size, (char *)response);
.588
           free(response):
.589
.591
```

... which brings us back to the terminal. Our call to yutani\_poll returns, giving us a KEY\_EVENT, which we then pass to the appropriately-named key\_event for further processing.

```
1121 void * handle incoming(void * garbage) {
1122
         while (!exit application) {
1123
             yutani msg t * m = yutani poll(yctx);
1124
             if (m) {
1125
                 switch (m->type) {
1126
                      case YUTANI_MSG_KEY_EVENT:
1127
1128
                              struct yutani msg key event * ke
1129
                              int ret = (ke->event.action == KE
1130
                              key event(ret, &ke->event);
1131
1132
                          break:
1133
                     case YUTANI_MSG_WINDOW_FOCUS_CHANGE:
1134
1135
                              struct yutani msg window focus ch
1136
                              yutani window t * win = hashmap q
1137
                              if (win) {
1138
                                  win->focused = wf->focused:
                                  render_decors();
1139
1140
1141
1142
                          break:
1143
                     case YUTANI MSG SESSION END:
1144
1145
                              kill(child pid, SIGKILL);
1146
                              exit application = 1;
1147
1148
                          break:
1149
                      case YUTANI MSG RESIZE OFFER:
1150
1151
                              struct yutani msg window resize *
1152
                              resize finish(wr->width, wr->heig
1153
1154
                          break:
1155
                     case YUTANI MSG WINDOW MOUSE EVENT:
```

key\_event checks if we need to do special processing before writing the key to the *TTY master*. Many keys, such as the *function keys and arrow keys*, need to be encoded as special *escape sequences* as they can't be represented as single bytes.

Our f key can just be sent as the character "f", though, so we call handle\_input('f') which will write to our TTY.



As there is now data available from our TTY, we can go back to the shell, which will finally return from fgetc, receiving the letter f.

30

32

34

37

38

42

43

44

45

47

48

54

59

62

To support advanced line editing like cursor movements, the shell needs to figure out if the data it is receiving is raw characters or escape sequences, so it runs everything through a state machine. Luckily, our f needs no special processing, so we can continue.

```
28      int rline(char * buffer, int buf size, rline callbacks t * call
       /* Initialize context */
       rline context t context = {
           buffer,
           callbacks,
           Θ,
           buf size,
           Θ,
           0,
           Θ.
       };
       printf("\033[s");
       fflush(stdout);
       key event state t kbd state = \{0\};
       /* Read keys */
       while ((context.collected < context.requested) && (!context</pre>
           uint32 t key sym = kbd key(&kbd state, fgetc(stdin));
           if (key sym == KEY NONE) continue;
           switch (key_sym) {
               case KEY CTRL C:
                   printf("^C\n");
                   context.buffer[0] = '\0';
                   return 0;
               case KEY CTRL R:
                    if (callbacks->rev search) {
                        callbacks->rev search(&context);
                        return context.collected;
                   continue;
               case KEY ARROW UP:
               case KEY CTRL P:
                   if (callbacks->key up) {
                        callbacks->kev un(&context):
```

Something we neglected to mention earlier, as we glossed over TTY modes, is that our TTY is set to not automatically echo the characters we type as we type them. This is important in supporting line editing, and means we need to write back out the character we received so it shows up on the terminal. We also need to add the `f` to an internal buffer.

```
} else {
    printf("%c", (char)key_sym);
    if (context.collected < context.requested) {
        context.buffer[context.collected] = (char)key_sym;
        context.buffer[++context.collected] = '\0';
        context.offset++;
    }
    fflush(stdout);</pre>
```

237

238

239

240

43

Since we've printed a character, we now need to go back up the TTY to the terminal again! Since we've already seen some of the details of the TTY device, we'll skip all the way to the terminal itself.

A thread in the terminal is always reading from the TTY master, and will receive our `f` and run it through the ANSI escape sequence processor.

```
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
```

```
unsigned char buf[1024];
while (!exit_application) {
    int r = read(fd_master, buf, 1024);
    spin_lock(&display_lock);
    for (uint32_t i = 0; i < r; ++i) {
        ansi_put(ansi_state, buf[i]);
    }
    display_flip();
    spin_unlock(&display_lock);
}
```

Terminals are complicated. In the early days of digital computer text, only basic black-and-white character cell displays were available, and they only supported a handful of characters.

These days, terminals support full-color Unicode text, hundreds of formatting options, cursor movement...

These formatting effects are made possible through complex escape sequences.



#### We're going to skip all of that.



At the end of the day, the escape processor will eventually call two functions: cell set and cell redraw. The first of these stores the `f` and all of its formatting hints into a data structure representing the terminal display, and the latter actually draws the glyph for `f` into the window for the terminal.

```
395 static void cell_set(uint16_t x, uint<mark>1</mark>6_t y, uint32_t c, uint32_t fg, uint32<u>t bg, uint8_t flags) {</u>
        if (x >= term_width || y >= term_height) return;
        term_cell_t * cell = (term_cell_t *)((uintptr_t)term_buffer + (y * term_width + x) * sizeof(term_cell_t));
        cell->c
                    = c:
                   = fg;
        cell->fa
        cell->bg
                  = bq;
        cell->flags = flags;
402
   static void cell redraw(uint16 t x, uint16 t y) {
        if (x >= term_width || y >= term_height) return;
        term cell t * cell = (term cell t *)((uintptr t)term buffer + (y * term width + x) * sizeof(term cell t));
407
        if (((uint32_t *)cell)[0] == 0x00000000) {
408
            term_write_char(' ', x * char_width, y * char_height, TERM_DEFAULT_FG, TERM_DEFAULT_BG, TERM_DEFAULT_FLAGS);
        } else {
            term_write_char(cell->c, x * char_width, y * char_height, cell->fg, cell->bg, cell->flags);
412
```

At this point, we need to take a quick break from looking at code and talk about how windows get displayed on the screen in ToaruOS.

As we mentioned earlier, ToaruOS has a *window compositor*. The most important job the compositor has is to *composite* - to take all of the windows from all of the running applications and turn them into a final scene.





ToaruOS's window compositor is called *Yutani* (a reference to Wayland and *Alien*). Yutani represents windows as *canvases* of 32-bit RGBA pixels, which it stores in shared memory regions.

When an application wants to draw to a window, it modifies the pixels directly in memory and then informs the compositor that it should redraw the modified regions using a *flip* message.





The compositor redraws only the regions which have been modified since the last redraw, producing a final image which it then copies into video memory.



To actually draw text, ToaruOS makes use of a library called Freetype to read TrueType and OpenType fonts and draw glyphs.

Once we've drawn the glyph for our 'f', we tell the compositor to flip our modified regions.

```
196 void drawChar(FT_Bitmap * bitmap, int x, int y, uint32_t fg, uint32_t bg) {
       int i, j, p, q;
       int x_max = x + bitmap->width;
199
       int y max = y + bitmap->rows;
       for (j = y, q = 0; j < y max; j++, q++) {
201
           for (i = x, p = 0; i < x max; i++, p++) {
               uint32 t a = ALP(fg);
203
               a = (a * bitmap->buffer[q * bitmap->width + p]) / 255;
204
               uint32_t tmp = rgba(_RED(fg),_GRE(fg),_BLU(fg),a);
               term_set_point(i,j, alpha_blend_rgba(premultiply(bg), premultiply(tmp)));
208 }
117 static void display flip(void) {
118
         if (l x != INT32 MAX && l y != INT32 MAX) {
             yutani_flip_region(yctx, window, l_x, l_y, r_x - l_x, r_y - l_y);
119
120
             l x = INT32 MAX;
             l y = INT32 MAX;
122
             r x = -1;
             r y = -1;
124
125 }
```

The compositor collects all of the screen updates and then *blits* the windows together.

```
1049 static void redraw windows(yutani globals t * yg) {
1050
         /* Save the cairo contexts so we can apply clipping */
1051
         save cairo states(yg);
1052
         int has updates = 0;
1053
1054
         /* We keep our own temporary mouse coordinates as they ma
1055
         int tmp mouse_x = yg->mouse_x;
1056
         int tmp mouse y = yq->mouse y;
1057
1058
         /* If the mouse has moved, that counts as two damage reg:
1059
         if ((yg->last_mouse_x != tmp_mouse_x) || (yg->last_mouse_
1060
             has updates = 2;
1061
             yutani_add_clip(yg, yg->last_mouse_x / MOUSE_SCALE -
1062
             yutani add clip(yg, tmp mouse x / MOUSE SCALE - MOUSE
1063
         }
1064
1065
         yq->last mouse x = tmp mouse x;
1066
         yg->last mouse y = tmp mouse y;
1067
1068
         if (yg->bottom z && yg->bottom z->anim mode) mark window
1069
         if (yg->top z && yg->top z->anim mode) mark window(yg, yc
1070
         foreach (node, yg->mid zs) {
1071
             yutani server window t * w = node->value;
1072
             if (w && w->anim mode) mark window(yg, w);
1073
         }
1074
1075
         /* Calculate damage regions from currently queued updates
1076
         spin lock(&yg->update list lock);
1077
         while (yg->update list->length) {
1078
             node t * win = list dequeue(yg->update list);
             yutani damage rect t * rect = (void *)win->value;
1079
1080
1081
             /* We add a clip region for each window in the update
1082
             has updates = 1;
1083
             yutani add clip(yg, rect->x, rect->y, rect->width, re
1084
             free(rect):
1085
             free(win):
1086
1087
         spin unlock(&yg->update list lock);
```

Using a library called *Cairo* to perform pixel pushing, the compositor can support transparent, rotated, and animated windows.

```
826 static int yutani blit window(yutani globals t * yg, cairo t * ctx,
827
828
        /* Obtain the previously initialized cairo contexts */
        cairo t * cr = ctx;
830
831
        /* Window stride is always 4 bytes per pixel... */
        int stride = window->width * 4:
        /* Initialize a cairo surface object for this window */
        cairo surface t * surf = cairo image surface create for data(
                window->buffer, CAIRO FORMAT ARGB32, window->width, winc
837
838
        /* Save cairo context */
839
        cairo save(cr);
840
841
842
         * Offset the rendering context appropriately for the position c
843
         * based on the modifier paramters
844
845
        cairo translate(cr, x, y);
846
847
        /* Top and bottom windows can not be rotated. */
848
        if (!window_is_top(yg, window) && !window_is_bottom(yg, window))
849
            /* Calcuate radians from degrees */
850
851
            /* XXX Window rotation is disabled until damage rects can ta
852
            if (window->rotation != 0) {
                double r = M_PI * (((double)window->rotation) / 180.0);
854
855
                /* Rotate the render context about the center of the wir
856
                cairo translate(cr, (int)( window->width / 2), (int)( (i
857
                cairo rotate(cr, r);
858
                cairo_translate(cr, (int)(-window->width / 2), (int)(-wi
859
860
                /* Prefer faster filter when rendering rotated windows *
861
                cairo pattern set filter (cairo get source (cr), CAIRO F
862
            }
863
864
        if (window->anim mode) {
865
            int frame = yutani time since(yg, window->anim start);
866
            if (frame >= yutani animation lengths[window->anim mode]) {
867
                  XXX handle animation-end things like cleanup of closi
868
                if (window->anim mode == YUTANI_EFFECT_FADE_OUT) {
```

The last step in this process is to copy the final frame to video memory so it shows up on screen.

* Flip the updated areas. This minimizes writes to video memory, * which is very important on real hardware where these writes are slow. */
<pre>cairo_set_operator(yg-&gt;real_ctx, CAIRO_OPERATOR_SOURCE); cairo_translate(yg-&gt;real_ctx, 0, 0);</pre>
cairo_set_source_surface[yg->real_ctx, yg->framebuffer_surface, 0, 0); cairo_paint(yg->real_ctx);
} spin unlock(&vg->redraw lock):

1151 1152

> 1154 1155

1156 1157

#### Finally, we've typed the letter `f`!



# This process repeats for the rest of our input until we hit enter.



At this point, the shell will parse our input, splitting it on spaces, expanding variables, etc.

```
526 int shell exec(char * buffer, int buffer size) {
527
528
        /* Read previous history entries */
        if (buffer[0] == '!') {
529
530
            uint32 t x = atoi((char *)((uintptr t)buffer + 1));
            if (x <= shell history_count) {</pre>
532
                buffer = shell history get(x - 1);
533
                buffer size = strlen(buffer);
534
            } else {
535
                fprintf(stderr, "esh: !%d: event not found\n", x);
536
                return 0;
537
            }
538
        }
539
540
        char * history = malloc(strlen(buffer) + 1);
541
        memcpy(history, buffer, strlen(buffer) + 1);
542
543
        if (buffer[0] != ' ' && buffer[0] != '\n') {
544
            shell history insert(history);
545
        } else {
546
            free(history);
547
        }
548
549
        char * argv[1024];
550
        int tokenid = 0:
```

## Once the shell has processed our input, it will fork a child process.

```
771  /* Now execute the last piece and wait on all of them */
772 } else {
773  shell_command_t func = shell_find(*arg_starts[0]);
774   if (func) {
775      return func(argcs[0], arg_starts[0]);
776      } else {
777           child_pid = fork();
778           if (!child_pid) {
779              run_cmd(arg_starts[0]);
780           }
781      }
782 }
```

Forking is the traditional Unix method for creating new processes. Forking produces a nearly identical copy of the running process - it has the same virtual memory contents, the same instruction pointer, the same file descriptors.

```
.82
183
     * @return To the parent: PID of the child; to the child: 0
186
187 uint32 t fork(void) {
        IRQ OFF;
188
189
190
        uintptr t esp, ebp;
191
        current_process->syscall_registers->eax = 0;
        /* Make a pointer to the parent process (us) on the stack */
        process t * parent = (process t *)current process;
        assert(parent && "Forked from nothing??");
        /* Clone the current process' page directory */
198
        page_directory_t * directory = clone_directory(current_directory);
199
        assert(directory && "Could not allocate a new page directory!");
        /* Spawn a new process from this one */
201
        debug print(INF0, "\033[1;32mALLOC {\033[0m");
202
203
        process_t * new_proc = spawn_process(current_process);
        debug print(INF0, "\033[1;32m}\033[0m");
204
        assert(new proc && "Could not allocate a new process!");
205
206
        /* Set the new process' page directory to clone */
        set process environment(new proc, directory);
207
        struct regs r;
209
        memcpy(&r, current_process->syscall_registers, sizeof(struct regs));
        new proc->syscall registers = &r;
211
212
        esp = new_proc->image.stack;
213
214
        ebp = esp;
215
        new proc->syscall registers->eax = 0;
216
217
        PUSH(esp, struct regs, r);
218
219
        new proc->thread.esp = esp;
220
        new proc->thread.ebp = ebp;
221
222
        new proc->is tasklet = parent->is tasklet;
```

But our new process will have a different return value for the call to fork than its parent, allowing it to take a different code path.

In the shell, that code path leads to a call to execvp. Within the C library, execvp is a complicated series of calls to readdin and stat to find a suitable executable within \$PATH. We'll skip that and jump to what happens when we call the underlying exec system call.

```
/* Now execute the last piece and wait on all of them */
771
772
       } else {
           shell command t func = shell find(*arg starts[0]);
774
           if (func) {
775
776
                return func(argcs[0], arg_starts[0]);
           } else {
                child pid = fork();
                if (!child pid) {
779
780
                    run cmd(arg_starts[0]);
                }
781
            }
782
        3
508 void run cmd(char ** args) {
        int i = execvp(*args, args);
509
510
        shell command t func = shell find(*args);
511
        if (func) {
            int argc = 0;
512
513
            while (args[argc]) {
514
                 argc++;
515
             i = func(argc, args);
516
517
        } else {
518
            if (i != 0) {
                 fprintf(stderr, "%s: Command not found\n", *args);
519
520
                 i = 127:
521
522
        }
523
        exit(i);
524 }
```

exec needs to load and parse our binary.

ToaruOS uses a binary format called ELF, and our first task is to read that ELF file and make sure it is actually an ELF.

So let's look at how the file system works.

```
16 int exec_elf(char * path, fs_node_t * file, int argc, char ** argv, char ** env) {
       Elf32_Header * header = (Elf32_Header *)malloc(file->length + 100);
       debug print(NOTICE, "---> Starting load.");
20
21
22
23
24
25
26
27
       IRQ RES;
       read fs(file, 0, file->length, (uint8 t *)header);
       IR0 OFF:
       debug print(NOTICE, "---> Finished load.");
       current process->name = malloc(strlen(path) + 1);
       memcpy(current process->name, path, strlen(path) + 1);
28
29
30
31
32
33
34
       current process->cmdline = argv;
       /* Alright, we've read the binary, time to load the loadable sections */
       /* Verify the magic */
       if (
               header->e_ident[0] != ELFMAG0 ||
                header->e_ident[1] != ELFMAG1
                header->e_ident[2] != ELFMAG2 |
                header->e ident[3] != ELFMAG3) {
36
37
           /* What? This isn't an ELF... */
           debug_print(ERROR, "Not a valid ELF executable.");
           free(header);
           close fs(file);
           return -1:
```

ToaruOS's primary on-disk filesystem is ext2, an older version of the most commonly used filesystem on Linux (ext4).

ext2 tracks files in a structure called an *inode*, and each *inode* references *blocks* that contain file data.

When we read part of a file, we need to find what blocks to look at and where they are on disk.



Our ext2 driver keeps recently used blocks in a cache, but if we haven't yet read the blocks we need, we'll need to read them from the underlying *block device* in this case, an ATA hard disk.

```
78 static uint32_t read_ata(fs_node_t *node, uint32_t offset, uint32_t size, uint8_t >
        struct ata_device * dev = (struct ata_device *)node->device;
       unsigned int start block = offset / ATA SECTOR SIZE;
       unsigned int end block = (offset + size - 1) / ATA SECTOR SIZE;
       unsigned int x offset = 0;
        if (offset > ata_max_offset(dev)) {
88
89
            return 0;
        if (offset + size > ata max offset(dev)) {
            unsigned int i = ata max offset(dev) - offset;
            size = i:
        3
 95
96
        if (offset % ATA_SECTOR_SIZE) {
            unsigned int prefix_size = (ATA_SECTOR_SIZE - (offset % ATA_SECTOR_SIZE));
            char * tmp = malloc(ATA_SECTOR_SIZE);
99
100
            ata device read sector(dev, start block, (uint8 t *)tmp);
101
            memcpy(buffer, (void *)((uintptr t)tmp + (offset % ATA SECTOR SIZE)), pref:
102
            free(tmp);
105
106
107
            x_offset += prefix_size;
            start_block++;
108
109
        if ((offset + size) % ATA SECTOR SIZE && start block < end block) {
110
            unsigned int postfix size = (offset + size) % ATA SECTOR SIZE;
111
            char * tmp = malloc(ATA SECTOR SIZE);
112
113
            ata_device_read_sector(dev, end_block, (uint8_t *)tmp);
            memcpy((void *)((uintptr_t)buffer + size - postfix_size), tmp, postfix_size
116
117
118
119
            free(tmp);
            end block--;
120
121
        while (start block <= end block) {</pre>
122
123
124
125
            ata_device_read_sector(dev, start_block, (uint8_t *)((uintptr_t)buffer + x_
            x_offset += ATA_SECTOR_SIZE;
            start block++;
126
```

Once we've read all of the blocks from disk, we can load our binary into memory by copying sections.

When we're done loading the binary, we jump to userspace and start running our newly loaded code.

```
/* Load the loadable segments from the binary */
       for (uintptr t x = 0; x < (uint32 t)header->e shentsize * header->e shnum; x += header->e sh
           /* read a section header */
           Elf32 Shdr * shdr = (Elf32 Shdr *)((uintptr t)header + (header->e shoff + x));
           if (shdr->sh addr) {
               /* If this is a loadable section, load it up. */
               if (shdr->sh addr == 0) continue; /* skip sections that try to load to 0 */
               if (shdr->sh_addr < current_process->image.entry) {
                   /* If this is the lowest entry point, store it for memory reasons */
                   current process->image.entry = shdr->sh addr;
               if (shdr->sh addr + shdr->sh size - current process->image.entry > current process->
60
                   /* We also store the total size of the memory region used by the application *,
                   current process->image.size = shdr->sh addr + shdr->sh size - current process->i
62
63
64
               for (uintptr t i = 0; i < shdr->sh size + 0x2000; i += 0x1000) {
                   /* This doesn't care if we already allocated this page */
                   alloc_frame(get_page(shdr->sh_addr + i, 1, current_directory), 0, 1);
66
                   invalidate tables at(shdr->sh addr + i);
67
               if (shdr->sh_type == SHT_NOBITS) {
                   /* This is the .bss, zero it */
                  memset((void *)(shdr->sh_addr), 0x0, shdr->sh_size);
               } else {
                   /* Copy the section into memory */
                  memcpy((void *)(shdr->sh_addr), (void *)((uintptr_t)header + shdr->sh_offset), s
       /* Store the entry point to the code comment */
155
           /* Go go go */
156
           enter_user_jmp(entry, argc, argv_, USER_STACK_TOP);
```

After all of that work, we're now running our fetch binary.

fetch makes use of a special feature in ToaruOS: The network filesystem. Instead of creating a socket, fetch opens the virtual file /dev/net/toaruos.org:80

```
6 int main(int argc, char * argv[]) {
      int opt:
      while ((opt = getopt(argc, argv, "?c:ho:")) != -1) {
          switch (opt) {
                   return usage(argv);
                  fetch options.cookie = optarg;
                  break;
              case 'h':
                  fetch options.show headers = 1;
                  break;
              case 'o':
                  fetch options.output file = optarg;
                  break:
94
      if (optind >= argc) {
          return usage(argv);
      struct http_req my_req;
      parse url(argv[optind], &my req);
      char file[100];
      sprintf(file, "/dev/net/%s", my_req.domain);
      fetch options.out = stdout;
      if (fetch options.output file) {
          fetch options.out = fopen(fetch options.output file, "w");
      FILE * f = fopen(file, "r+");
      if (!f) {
          fprintf(stderr, "Nope.\n");
          return 1:
      if (fetch_options.cookie) {
          fprintf(f,
               "GET /%s HTTP/1.0\r\n"
               "User-Agent: curl/7.35.0\r\n"
               "Host: %s\r\n'
               "Accept: */*\r\n"
               "Cookie: %s\r\n"
               "\r\n", my req.path, my req.domain, fetch options.cookie);
      } else {
          fprintf(f,
               "GET /%s HTTP/1.0\r\n"
               "User-Agent: curl/7.35.0\r\n"
               "Host: %s\r\n"
               "Accept: */*\r\n"
               "\r\n", my req.path, my req.domain);
```

In the network filesystem, we lookup "toaruos.org" and create a TCP socket. We store this socket in a *virtual file node* that we'll add to the running process's file descriptor table.

```
291 static fs_node_t * finddir_netfs(fs_node_t * node, char * name) {
        debug_print(WARNING, "Need to look up domain or check if is IP: %s", nam
296
297
        int port = 80;
        char * colon;
        if ((colon = strstr(name, ":"))) {
             *colon = '\0';
301
302
303
304
305
306
             colon++:
             port = atoi(colon);
        uint32 t ip = 0;
        if (is ip(name)) {
            debug print(WARNING, "
                                       IP: %x", ip_aton(name));
308
309
310
            ip = ip_aton(name);
        } else {
             if (hashmap has(dns cache, name)) {
                 ip = ip_aton(hashmap_get(dns_cache, name));
                 debug print(WARNING, " In Cache: %s → %x", name, ip);
            } else {
                 debug print(WARNING, "
                                            Still needs look up.");
315
316
                 return NULL:
        fs node t * fnode = malloc(sizeof(fs node t));
320
321
322
323
324
325
326
327
328
329
330
331
332
        memset(fnode, 0x00, sizeof(fs node t));
        fnode->inode = 0:
        strcpy(fnode->name, name);
        fnode->mask = 0555;
        fnode->flags = FS CHARDEVICE;
                         = socket read;
        fnode->read
        fnode->write
                       = socket write;
        fnode->device = (void *)net_open(SOCK_STREAM);
        net connect((struct socket *)fnode->device, ip, port);
        return fnode:
```

#### Writing to the virtual file that the network filesystem gave us will send TCP packets.



## Back in fetch, we write our HTTP requests headers...

```
76 int main(int argc, char * argv[]) {
         int opt;
         while ((opt = getopt(argc, argv, "?c:ho:")) != -1) {
             switch (opt) {
                       return usage(argv);
                       fetch_options.cookie = optarg;
                       break;
                       fetch options.show headers = 1;
                       break;
                       fetch options.output file = optarg;
 92
93
                       break:
         if (optind >= argc) {
             return usage(argv);
 98
99
100
102
103
104
105
106
        struct http_req my_req;
parse_url(argv[optind], &my_req);
         char file[100];
         sprintf(file, "/dev/net/%s", my_req.domain);
         fetch options.out = stdout;
        if (fetch_options.output_file) {
    fetch_options.output_file) {
    fetch_options.out = fopen(fetch_options.output_file, "w");
}
         }
110
111
         FILE * f = fopen(file, "r+");
         if (!f) {
             fprintf(stderr, "Nope.\n");
             return 1:
         if (fetch_options.cookie) {
             fprintf(f,
                   "GET /%s HTTP/1.0\r\n"
                  "User-Agent: curl/7.35.0\r\n"
                  "Host: %s\r\n"
                  "Accept: */*\r\n"
                  "Cookie: %s\r\n"
                  "\r\n", my_req.path, my_req.domain, fetch_options.cookie);
        } else {
128
129
             fprintf(f,
                   "GET /%s HTTP/1.0\r\n"
                  "User-Agent: curl/7.35.0\r\n"
                  "Host: %s\r\n"
                  "Accept: */*\r\n"
                  "\r\n", my req.path, my req.domain);
```

When the network card receives packets from the router, it passes them to the *network* processing thread. Each open socket has a corresponding queue of data that has yet to be read, and the network processing thread will add our newly received packet to the appropriate queue.

Reading from the virtual file will call net\_recv to pull data out of the queue.

```
515 size_t net_recv(struct socket* socket, uint8_t* buffer, size t len) {
        tcpdata t *tcpdata = NULL;
        node t *node = NULL;
        debug print(WARNING, "0x%x [socket]", socket);
521
        size t offset = 0;
        size t size to read = 0;
523
524
525
526
527
528
529
530
        do {
        if (socket->bytes available) {
            tcpdata = socket->current packet;
        } else {
            spin lock(socket->packet queue lock);
531
532
533
                 if (socket->packet queue->length > 0) {
                     node = list dequeue(socket->packet queue);
                     spin unlock(socket->packet queue lock);
 34
35
36
                     break;
                } else {
                     if (socket->status == 1) {
                         spin unlock(socket->packet queue lock);
                         debug print(WARNING, "Socket closed, done reading.");
539
540
541
542
                         return 0:
                     spin unlock(socket->packet queue lock);
                     sleep on(socket->packet wait);
                     spin lock(socket->packet queue lock);
 544
 45
            } while (1);
            tcpdata = node->value;
            socket->bytes_available = tcpdata->payload size;
            socket->bytes read = 0;
 50
            free(node);
 52
53
54
        size to read = MIN(len, offset + socket->bytes available);
 555
556
557
        if (tcpdata->payload != 0) {
            memcpy(buffer + offset, tcpdata->payload + socket->bytes read, size to read);
        offset += size to read;
        if (size to read < socket->bytes available) {
            socket->bytes_available = socket->bytes_available - size_to_read;
            socket->bytes read = size to read;
            socket->current packet = tcpdata:
```

#### The web server on the other end will respond, and we will read the result and print it to the terminal.

```
http_parser_settings settings;
136
137
        memset(&settings, 0, sizeof(settings));
138
        settings.on_header_field = callback_header_field;
139
        settings.on header value = callback header value;
140
        settings.on_body = callback_body;
141
142
       http_parser parser;
143
        http_parser_init(&parser, HTTP_RESPONSE);
144
145
       while (!feof(f)) {
146
            char buf[1024];
47
            memset(buf, 0, sizeof(buf));
148
            size_t r = fread(buf, 1, 1024, f);
149
            http parser execute(&parser, &settings, buf, r);
150
        }
151
152
```

```
fflush(fetch_options.out);
```

Turns out we just printed a PDF to our terminal, which looks like a bunch of garbage.

Oops.





#### Demos

Python 2.7

vim, gcc

Quake



## Design Goals - Now

Porting software

GUI

Clean code

#### Future

Now: "Misaka" x86-64 kernel project

Next year: Porting NetSurf; maybe Webkit?

Far future: glib, GTK...



#### Thanks!

**IRC** #toaruos on Freenode

Sites toaruos.org github.com/klange/toaruos

Historical Screenshots y/toaruos-screens

**Live CD** y/toaruos-iso

