# What Tilck is

- A project consisting on:
  - A monolithic kernel written in C and assembly
  - A bootloader working both on UEFI and legacy BIOS systems
  - Several test suites and a powerful CMake-based build system
  - Buildroot-like scripts for downloading & building 3rd party software
- Partially compatible with Linux at binary level
- Uniprocessor, but fully preemptable
- Educational, with potential to be more than that (see testing etc.)
- Runs only on i686, at the moment (will be ported to ARM, RISC-V etc.)
- Open source, distributed under the BSD 2-clause license

# What Tilck is **NOT**

▶ An attempt to replace Linux

▶ An attempt to be yet another desktop operating system

▶ An attempt to be a large-scale server operating system

▶ A real-time OS, but it might become one in the future

▶ A OS running on NOMMU machines, but (probably) will in the future

▶ Ready for production use: it still lacks features as storage, networking etc.

# Why the binary compatibility with Linux?

▶ It's cool being able to test the same "bits" both on Linux and Tilck

▶ Robustness: Tilck can empirically show robustness and correctness by running 3rd party software never written for it

▶ Didn't want to design a whole new syscall interface from scratch

▶ Didn't want to implement a whole libc too

▶ Didn't want to build a custom GCC toolchain. I wanted to use the pre-built toolchains from: https://toolchains.bootlin.com/

▶ Increase the likelihood the project to get more interest from the community?

▶ Porting pre-existing software to Tilck will require a little or no effort at all.

# Core values & goals

▶ Minimal memory footprint

▶ Ultra low-latency

▶ Deterministic behavior

▶ Extra robustness

▶ Portability

▶ Simplicity

▶ Partial compatibility with Linux

▶ Must work on real (modern) hardware

▶ Exceptional developer experience: building & testing the project should be as easy as technologically possible

# Live demo

Because a demo is worth more than a thousand words

Vladislav K. Valtchev (2022)

# Funny stories & interesting challenges

Vladislav K. Valtchev (2022)

# My latest bug [1/6]

*(and its 2-char fix)*

▶ I have a test (fork_oom) that:

1. Estimates the amount of *committed* memory that can be used
2. Allocates and commits more than half of that
3. Calls fork()
4. In the child, tries to commit *all* of that memory
5. Expects the child to be killed by the kernel

# My latest bug [1/6]

*(and its 2-char fix)*

▶ I have a test (fork_oom) that:

1. Estimates the amount of *committed* memory that can be used

2. Allocates and commits more than half of that

3. Calls fork()

4. In the child, tries to commit *all* of that memory

5. Expects the child to be killed by the kernel

▶ I just found that it fails on real HW machines

# My latest bug [1/6]

*(and its 2-char fix)*

▶ I have a test (fork_oom) that:

1. Estimates the amount of *committed* memory that can be used

2. Allocates and commits more than half of that

3. Calls fork()

4. In the child, tries to commit *all* of that memory

5. Expects the child to be killed by the kernel

▶ I just found that it fails on real HW machines

▶ Quickly, I discovered that it fails on VMs too but only when they have significantly more RAM. That's weird. Mmm...

# My latest bug [1/6]

*(and its 2-char fix)*

▶ I have a test (fork_oom) that:

1. Estimates the amount of *committed* memory that can be used

2. Allocates and commits more than half of that

3. Calls fork()

4. In the child, tries to commit *all* of that memory

5. Expects the child to be killed by the kernel

▶ I just found that it fails on real HW machines

▶ Quickly, I discovered that it fails on VMs too but only when they have significantly more RAM. That's weird. Mmm…

▶ I had to debug that.

Vladislav K. Valtchev (2022)

# My latest bug [2/6]



That's fine...

How could we commit so much memory?
262 MB x 2 = 524 MB > 501 MB [usable]
(ehm, we don't have swap)

That means trying to free a page not allocated in the heap, during munmap().

# My latest bug [3/6]

*So, I started debugging the CoW page-fault logic…*

```c
2 bool handle_potential_cow(void *context)
3 {
4     /* ... */
5
6     if (!(pt->pages[pt_index].avail & PAGE_COW_ORIG_RW))
7         return false; /* Not a COW page */
8
9     const u32 orig_page_paddr = (u32)
10        pt->pages[pt_index].pageAddr << PAGE_SHIFT;
11
12    if (pf_ref_count_get(orig_page_paddr) == 1) {
13
14        /* This page is not shared anymore. No need for copying it. */
15        pt->pages[pt_index].rw = true;
16        pt->pages[pt_index].avail = 0;
17        invalidate_page_hw(vaddr);
18        return true;
19    }
20
21    /* ... */
22 }
```

After committing a few MBs in the child, we end up here!

# My latest bug [4/6]

*I realized I had ASSERTs disabled in that build! So, after turning them on...*



Aha, gotcha! You're really trying to free the **zero page**!

# My latest bug [5/6]

*Let's look at this limit case…*



```
root@tilck:/# devshell -c fork_oom
[devshell] Executing built-in command 'fork_oom'
[   parent   ] Estimating usable memory..
[child] Pid: 37
[   7.798] Out-of-memory: killing pid 37
[   parent   ] Child killed by signal 9
[   parent   ] Estimated usable memory: 487 MB
Alloc 255 MB...
Write to the buffer...
Done. Now, fork()..
Child [38]: writing to the whole CoW buffer...
[  10.270] Out-of-memory: killing pid 38
parent: the child exited with signal 9, as expected.
root@tilck:/#
```

Allocating 255 MB works…

# My latest bug [6/6]

*That means only one thing...*

```
 2  static u16 *pageframes_refcount;
 3  static ulong phys_mem_lim;
 4
 5  static ALWAYS_INLINE u32 pf_ref_count_get(u32 paddr)
 6  {
 7      if (UNLIKELY(paddr >= phys_mem_lim))
 8          return 0;
 9
10      return pageframes_refcount[paddr >> PAGE_SHIFT];
11  }
```

# My latest bug [6/6]

*That means only one thing...*

That's the problem: a 16-bit ref-count

```
 2  static u16 *pageframes_refcount;
 3  static ulong phys_mem_lim;
 4
 5  static ALWAYS_INLINE u32 pf_ref_count_get(u32 paddr)
 6  {
 7      if (UNLIKELY(paddr >= phys_mem_lim))
 8          return 0;
 9
10      return pageframes_refcount[paddr >> PAGE_SHIFT];
11  }
```

# My latest bug [6/6]

*That means only one thing...*

That's the problem: a 16-bit ref-count

```
2  static u16 *pageframes_refcount;
3  static ulong phys_mem_lim;
4
5  static ALWAYS_INLINE u32 pf_ref_count_get(u32 paddr)
6  {
7     if (UNLIKELY(paddr >= phys_mem_lim))
8        return 0;
9
10    return pageframes_refcount[paddr >> PAGE_SHIFT];
11 }
```

It wraps around after 65,535 pages, meaning that the kernel cannot support 256 MB or more of *uncommited* memory!

# Making the framebuffer console *fast*

Vladislav K. Valtchev (2022)

# Making the framebuffer console *fast*

- Premise: why implement a framebuffer console?
  - Text mode was almost completely dead even 5 years ago
  - Pure-UEFI machines don't support text mode
  - Text mode is a x86 thing: Raspberry PI and other machines don't support it

# Making the framebuffer console *fast*

- Premise: why implement a framebuffer console?
  - Text mode was almost completely dead even 5 years ago
  - Pure-UEFI machines don't support text mode
  - Text mode is a x86 thing: Raspberry PI and other machines don't support it

- Why speed matters so much? Just mark the pages as WC and it will be reasonably fast.

# Making the framebuffer console *fast*

- Premise: why implement a framebuffer console?
  - Text mode was almost completely dead even 5 years ago
  - Pure-UEFI machines don't support text mode
  - Text mode is a x86 thing: Raspberry PI and other machines don't support it

- Why speed matters so much? Just mark the pages as WC and it will be reasonably fast.
  - I didn't know about WC (write-combining) at the time

# Making the framebuffer console *fast*

- ► Premise: why implement a framebuffer console?
  - ► Text mode was almost completely dead even 5 years ago
  - ► Pure-UEFI machines don't support text mode
  - ► Text mode is a x86 thing: Raspberry PI and other machines don't support it

- ► Why speed matters so much? Just mark the pages as WC and it will be reasonably fast.
  - ► I didn't know about WC (write-combining) at the time
  - ► Therefore, I implemented a series of optimizations before discovering WC

# PSF fonts: a bitfield per each glyph



8 bit

16 bit

8 bit  8 bit

32 bit

# The simplest draw function (failsafe)

```c
4   static inline void fb_draw_pixel(u32 x, u32 y, u32 color)
5   {
6       if (fb_bpp == 32)
7           *(volatile u32 *) (fb_vaddr + (fb_pitch * y) + (x << 2)) = color;
8       else
9           // Assumption: bpp is 24
10          memcpy((void *) (fb_vaddr + (fb_pitch * y) + (x * 3)), &color, 3);
11  }
12
13  void fb_draw_char(u32 x, u32 y, u16 e)
14  {
15      u8 *data = font_glyph_data + font_bytes_per_glyph * vgaentry_get_char(e);
16      u32 arr[] = { vga_rgb_colors[vgaentry_get_fg(e)], vga_rgb_colors[vgaentry_get_bg(e)] };
17
18      for (u32 row = y; row < (y + font_h); row++, data += font_width_bytes) {
19          for (u32 b = 0; b < font_width_bytes; b++) {
20              for (u32 i = 0; i < 8; i++)
21                  fb_draw_pixel(x + (b << 3) + (8 - i - 1),   /* x */
22                                row,                           /* y */
23                                arr[!(data[b] & (1 << i))]);   /* color */
24          }
25      }
26  }
```

# Performance? Too slow, in particular on the modern machine (left)

16x8 font, 800x600

## Intel Core i7-7500U Kaby Lake

▶ 1,124,773 RDTSC cycles / char (avg.) [~385.7 µs]

## Intel Atom N270 Diamondville (32-bit)

▶ 297,287 RDTSC cycles / char (avg.) [~186.3 µs]

32x16 font, 3200x1800

▶ 7,416,012 RDTSC cycles / char (avg.) [~2543.2 µs]

Scrolling the whole screen takes several seconds!!

# A naïve optimization: loop unrolling

```
 3 #define draw_char_partial(b)                                                 \
 4    do {                                                                        \
 5       fb_draw_pixel(x + (b << 3) + 7, row, arr[!(data[b] & (1 << 0))]);       \
 6       fb_draw_pixel(x + (b << 3) + 6, row, arr[!(data[b] & (1 << 1))]);       \
 7       fb_draw_pixel(x + (b << 3) + 5, row, arr[!(data[b] & (1 << 2))]);       \
 8       fb_draw_pixel(x + (b << 3) + 4, row, arr[!(data[b] & (1 << 3))]);       \
 9       fb_draw_pixel(x + (b << 3) + 3, row, arr[!(data[b] & (1 << 4))]);       \
10       fb_draw_pixel(x + (b << 3) + 2, row, arr[!(data[b] & (1 << 5))]);       \
11       fb_draw_pixel(x + (b << 3) + 1, row, arr[!(data[b] & (1 << 6))]);       \
12       fb_draw_pixel(x + (b << 3) + 0, row, arr[!(data[b] & (1 << 7))]);       \
13    } while (0)
14
15 void fb_draw_char(u32 x, u32 y, u16 e)
16 {
17    u8 *data = font_glyph_data + font_bytes_per_glyph * vgaentry_get_char(e);
18    u32 arr[] = { vga_rgb_colors[vgaentry_get_fg(e)], vga_rgb_colors[vgaentry_get_bg(e)] };
19
20    if (LIKELY(font_width_bytes == 1))
21       for (u32 row = y; row < (y+font_h); row++, data += font_width_bytes)
22          draw_char_partial(0);
23    else if (font_width_bytes == 2)
24       for (u32 row = y; row < (y+font_h); row++, data += font_width_bytes) {
25          draw_char_partial(0);
26          draw_char_partial(1);
27       }
28    else
29       for (u32 row = y; row < (y+font_h); row++, data += font_width_bytes)
30          for (u32 b = 0; b < font_width_bytes; b++)
31             draw_char_partial(b);
32 }
```

# Benefits? Nah.

## Intel Core i7-7500U Kaby Lake

| Before (avg.) | 385.72 μs / char |
|---|---|
| After (avg.) | 384.44 μs / char |
| Speed up | 0.3% faster |

## Intel Atom N270 Diamondville (32-bit)

| Before (avg.) | 186.27 μs / char |
|---|---|
| After (avg.) | 175.30 μs / char |
| Speed up | 6.2% faster |

> Old school optimizations work better on old school machines!

# Intuition 1: rendering glyphs pixel by pixel is too slow

Vladislav K. Valtchev (2022)

# Solution 1: pre-rendering!

▶ But… is pre-rendering *every* glyph in the font even feasible?

framebuffer

# Pre-rendering! (font 16x8)

## 16 x 8 x 4 x 256 x 16 x 16 =

Height x Width

Bytes per pixel

# glyphs

FG colors

BG colors

## 32 MB: unfeasible!

# Pre-rendering! (font 32x16)

32 x 16 x 4 x 256 x 16 x 16 =

Height x Width

Bytes per pixel

\# glyphs

FG colors

BG colors

## 128 MB: pure madness!

# A better idea: pre-render all the possible 8-bit "scanlines" (= glyph rows)

$$2^8 \text{ x } 4 \text{ x } 8 \text{ x } 16 \text{ x } 16 = \text{ 2 MB}$$

All scanlines

Bytes per pixel

Scanline length

FG colors

BG colors

Still expensive, but affordable!

# It works on 32x16 fonts too!

# The pre-render code

```c
3  #define PSZ          4      /* pixel size = 32 bpp / 8 = 4 bytes */
4  #define SL_COUNT   256      /* all possible 8-pixel scanlines */
5  #define SL_SIZE      8      /* scanline size: 8 pixels */
6  #define FG_COLORS   16      /* #fg colors */
7  #define BG_COLORS   16      /* #bg colors */
8  #define TOT_CHAR_SCANLINES_SIZE (PSZ*SL_COUNT*FG_COLORS*BG_COLORS*SL_SIZE)

10 bool fb_pre_render_char_scanlines(void)
11 {
12     fb_w8_char_scanlines = kmalloc(TOT_CHAR_SCANLINES_SIZE);

14     if (!fb_w8_char_scanlines)
15         return false;

17     for (u32 fg = 0; fg < FG_COLORS; fg++) {
18         for (u32 bg = 0; bg < BG_COLORS; bg++) {
19             for (u32 sl = 0; sl < SL_COUNT; sl++) {
20                 for (u32 pix = 0; pix < SL_SIZE; pix++) {
21                     fb_w8_char_scanlines[
22                         fg * (BG_COLORS * SL_COUNT * SL_SIZE) +
23                         bg * (SL_COUNT * SL_SIZE) +
24                         sl * SL_SIZE +
25                         (SL_SIZE - pix - 1)
26                     ] = (sl & (1 << pix)) ? vga_rgb_colors[fg] : vga_rgb_colors[bg];
27                 }
28             }
29         }
30     }
31     return true;
32 }
```

# Intuition 2: copying 4 bytes at a time is too slow!

▶ Pre-rendering the glyphs or the just the "scanlines" is **not enough**

▶ The x86 `rep movsl` instruction copies just 4 bytes (= 1 pixel) at a time

# Solution 2: use the FPU

▶ Introduce something like **fpu_memcpy()**

▶ Write a whole row at a time during scrolling

▶ Only this way, we could offset the cost of saving/restoring the FPU registers

```
1  void fb_draw_row(u32 y, u16 *entries, u32 count, bool fpu)
2  {
3      static const void *ops[] = {
4          &&width_1_nofpu, &&width_1_fpu, &&width_2_nofpu, &&width_2_fpu
5      };
6
7      const u32 bpg_shift = 4 + (font_bytes_per_glyph == 64) * 2; // 4 or 6
8      const u32 w4_shift  = 5 + (font_w == 16);                   // 5 or 6
9      const void *const op = ops[(font_w == 16) * 2 + fpu];       // ops[0..3]
10     const ulong vaddr_base = fb_vaddr + (fb_pitch * y);
11
12     for (u32 ei = 0; ei < count; ei++) {
13
14         const u16 e = entries[ei];
15         const u32 c_off = (u32) ((vgaentry_get_fg(e) << 15) + (vgaentry_get_bg(e) << 11));
16         void *vaddr = (void *)vaddr_base + (ei << w4_shift);
17         const u8 *d = &font_glyph_data[vgaentry_get_char(e) << bpg_shift];
18         u32 *scanlines = &fb_w8_char_scanlines[c_off];
19         goto *op;
20
21     width_1_fpu:
22         for (u32 r = 0; r < font_h; r++, d++, vaddr += fb_pitch)
23             fpu_cpy_single_256_nt(vaddr, &scanlines[d[0] << 3]);
24         continue;
25
26     width_1_nofpu:
27         for (u32 r = 0; r < font_h; r++, d++, vaddr += fb_pitch)
28             memcpy32(vaddr, &scanlines[d[0] << 3], SL_SIZE);
29         continue;
30
31     width_2_fpu:
32         for (u32 r = 0; r < font_h; r++, d+=2, vaddr += fb_pitch) {
33             fpu_cpy_single_256_nt(vaddr,      &scanlines[d[0] << 3]);
34             fpu_cpy_single_256_nt(vaddr + 32, &scanlines[d[1] << 3]);
35         }
36         continue;
37
38     width_2_nofpu:
39         for (u32 r = 0; r < font_h; r++, d+=2, vaddr += fb_pitch) {
40             memcpy32(vaddr,      &scanlines[d[0] << 3], SL_SIZE);
41             memcpy32(vaddr + 32, &scanlines[d[1] << 3], SL_SIZE);
42         }
43         continue;
44     }
45 }
```

Flag: during IRQ, we cannot use the FPU

Scanlines for the given FG/BG colors

Jump to the same address during the whole loop

Copy 256 bit (32 bytes) the fastest way possible

# The FPU code [1/2]

```
3  static void *get_fpu_cpy_single_256_nt_func(void)
4  {
5      if (!kopt_no_fpu_memcpy) {
6
7          if (x86_cpu_features.can_use_avx2)
8              return &fpu_cpy_single_256_nt_avx2;
9
10         if (x86_cpu_features.can_use_sse2)
11             return &fpu_cpy_single_256_nt_sse2;
12
13         if (x86_cpu_features.can_use_sse)
14             return &fpu_cpy_single_256_nt_sse;
15     }
16
17     return NULL;
18 }
19
20 void init_fpu_memcpy(void)
21 {
22     void *func;
23
24     if ((func = get_fpu_cpy_single_256_nt_func())) {
25         simple_hot_patch(&__asm_fpu_cpy_single_256_nt, func, 128);
26     }
27 }
```

# The FPU code [2/2]

```
30 ALWAYS_INLINE FASTCALL void
31 fpu_cpy_single_256_nt_avx2(void *dest, const void *src)
32 {
33     asmVolatile("vmovdqa    (%0), %%ymm0\n\t"
34                 "vmovntdq %%ymm0,    (%1)\n\t"
35                 : /* no output */
36                 : "r" (src), "r" (dest)
37                 : "memory");
38 }
39
40
41 ALWAYS_INLINE FASTCALL void
42 fpu_cpy_single_256_nt_sse2(void *dest, const void *src)
43 {
44     asmVolatile("movdqa    (%0), %%xmm0\n\t"
45                 "movdqa 16(%0), %%xmm1\n\t"
46                 "movntdq %%xmm0,    (%1)\n\t"
47                 "movntdq %%xmm1, 16(%1)\n\t"
48                 : /* no output */
49                 : "r" (src), "r" (dest)
50                 : "memory");
51 }
```

```
53 ALWAYS_INLINE FASTCALL void
54 fpu_cpy_single_256_nt_sse(void *dest, const void *src)
55 {
56     asmVolatile("movq (%0), %%mm0\n\t"
57                 "movq 8(%0), %%mm1\n\t"
58                 "movq 16(%0), %%mm2\n\t"
59                 "movq 24(%0), %%mm3\n\t"
60                 "movntq %%mm0, (%1)\n\t"
61                 "movntq %%mm1, 8(%1)\n\t"
62                 "movntq %%mm2, 16(%1)\n\t"
63                 "movntq %%mm3, 24(%1)\n\t"
64                 : /* no output */
65                 : "r" (src), "r" (dest)
66                 : "memory");
67 }
```

# The moment of truth

*Font 16x8, resolution 800x600, default memory type\*, not WC*

*\* Typically that means UC (uncacheable) set through MTRRs*

Core i7-7500U Kaby Lake, AVX 2, 256-bit regs

| Before (avg.) | 385.72 µs / char |
|---|---|
| After (avg.) | 67.42 µs / char |
| Speed up | 5.72x faster |

Atom N270 Diamondville (32-bit), SSSE 3, 128-bit regs

| Before (avg.) | 186.27 µs / char |
|---|---|
| After (avg.) | 94.82 µs / char |
| Speed up | 1.96x faster |

Smaller impact, but smaller regs here

## Not bad at all!

# The moment of truth

*Font 16x8, resolution 800x600, default memory type\*, not WC*

*\* Typically that means UC (uncacheable) set through MTRRs*

Core i7-7500U Kaby Lake, AVX 2, 256-bit regs

| | |
|---|---|
| Before (avg.) | 385.72 µs / char |
| After (avg.) | 67.42 µs / char |
| Speed up | 5.72x faster |

Atom N270 Diamondville (32-bit), SSSE 3, 128-bit regs

| | |
|---|---|
| Before (avg.) | 186.27 µs / char |
| After (avg.) | 94.82 µs / char |
| Speed up | 1.96x faster |

*Font 32x16, resolution 3200x1800, default memory type\*, not WC*

| | |
|---|---|
| Before (avg.) | 2543.21 µs / char |
| After (avg.) | 371.54 µs / char  ⟵ *Still not fast enough, though* |
| Speed up | 6.84x faster |

Wow, that's close to the max 8x improvement!
(From 32 bit/write to 256 bit/write)

# The writing combining memory type (WC)

▶ Allows data to combined, temporarily stored in a buffer (WCB) and then released in burst mode

▶ Cannot be used most of the time because offers just **weak ordering**

▶ Can be set using **PAT** or **MTRR**s

▶ It's perfect for frame buffers

# Performance: the full picture [modern machine]
*Font 16x8, resolution 800x600, 32 bbp*

Intel Core i7-7500U Kaby Lake (AVX 2, 256-bit fpu regs)

| Mode | Opt | Wc | FPU | Cycles/char | usec/char |
|---|---|---|---|---|---|
| **Failsafe slow** | | | | 1,124,773 | 385.72 |
| **Failsafe opt** | | | | 1,121,034 | 384.44 |
| **Opt + fpu** | ✓ | | ✓ | 196,584 | 67.42 |
| Opt (no fpu) | ✓ | | | 1,177,902 | 403.94 |
| **Wc** | | ✓ | | 34,055 | 11.68 |
| Opt + Wc (no fpu) | ✓ | ✓ | | 44,294 | 15.19 |
| **Opt + Wc + fpu** | ✓ | ✓ | ✓ | 30,271 | 10.38 |

**32.9x faster!**

Just 12.5% faster

Vladislav K. Valtchev (2022)

# Performance: the full picture [older machine]
*Font 16x8, resolution 800x600, 32 bbp*

Intel Atom N270 Diamondville (32-bit, SSSE 3, 128 bit fpu regs)

| Mode | Opt | Wc | FPU | Cycles/char | usec/char |
|---|---|---|---|---|---|
| **Failsafe slow** | | | | 297,287 | 186.27 |
| **Failsafe opt** | | | | 279,781 | 175.30 |
| **Opt + fpu** | ✓ | | ✓ | 151,337 | 94.82 |
| Opt (no fpu) | ✓ | | | 279,453 | 175.10 |
| **Wc** | | ✓ | | 136,925 | 85.79 |
| Opt + Wc (no fpu) | ✓ | ✓ | | 136,914 | 85.79 |
| **Opt + Wc + fpu** | ✓ | ✓ | ✓ | 136,906 | 85.78 |

**2.04x faster**

No difference at all!

Vladislav K. Valtchev (2022)

# Performance on native res [modern machine]
*Font 32x16, resolution 3200x1800, 32 bbp*

Intel Core i7-7500U Kaby Lake (AVX 2, 256-bit fpu regs)

| Mode | Opt | Wc | FPU | Cycles/char | usec/char |
|---|---|---|---|---|---|
| **Failsafe slow** | | | | 7,416,012 | 2543.21 |
| **Failsafe opt** | | | | 7,408,637 | 2540.68 |
| **Opt + fpu** | ✓ | | ✓ | 1,083,408 | 371.54 |
| Opt (no fpu) | ✓ | | | 7,815,696 | 2680.28 |
| **Wc** | | ✓ | | 73,165 | 25.09 |
| Opt + Wc (no fpu) | ✓ | ✓ | | 159,067 | 54.55 |
| **Opt + Wc + fpu** | ✓ | ✓ | ✓ | 27,841 | 9.55 |

6.84x faster

**101.26x faster!**

2.63x faster

Not bad!

# Performance vs Linux [modern machine]

*Font 32x16, resolution 3200x1800, 32 bbp*



*Commit a858f229, release build*

*Kernel 5.4.0 (Ubuntu 20.04.4 LTS)*

▶ 9.55 μs / char

# Performance vs Linux [modern machine]
*Font 32x16, resolution 3200x1800, 32 bbp*



*Commit a858f229, release build*



*Kernel 5.4.0 (Ubuntu 20.04.4 LTS)*

▶ 9.55 μs / char

▶ 56.40 μs / char

# Performance vs Linux [modern machine]
*Font 32x16, resolution 3200x1800, 32 bbp*



*Commit a858f229, release build*

▶ 9.55 μs / char

*Kernel 5.4.0 (Ubuntu 20.04 LTS)*

▶ 56.40 μs / char

**5.9x faster!**

# Performance vs Linux [modern machine]

*Font 32x16, resolution 3200x1800, 32 bbp*



*Commit a858f229, release build*

*Kernel 5.4.0 (Ubuntu 20.04 LTS)*

56.40 μs – Linux 5.4.0

**?**

25.09 μs – Tilck failsafe + WC

9.55 μs – Tilck's best OPT

Vladislav K. Valtchev (2022)

## The benchmark code

```c
void console_perf_test(void)
{
    static const char letters[] =
        "ABCDEFGHIJKLMNOPQRSTUVWXYZ"
        "abcdefghijklmnopqrstuvwxyz0123456789";

    int iters = 3;
    struct winsize w;
    char *buf, tot_time_s[32], c_time_s[32];
    ssize_t r, tot, written;
    struct timespec ts_before, ts_after;
    uint64_t start, end, c;
    double tot_time_real, tot_time, time_c, cycles_per_sec;

    if (ioctl(1, TIOCGWINSZ, &w) != 0) {
        perror("ioctl() failed");
        return;
    }

    tot = w.ws_row * w.ws_col;

    if (!(buf = malloc(tot))) {
        perror("malloc() failed\n");
        return;
    }

    for (int i = 0; i < tot; i++) {
        buf[i] = letters[i % (sizeof(letters) - 1)];
    }

    printf("%s", CSI_ERASE_DISPLAY CSI_MOVE_CURSOR_TOP_LEFT);

retry:
    clock_gettime(CLOCK_REALTIME, &ts_before);
    start = RDTSC();

    for (int i = 0; i < iters; i++) {
        for (r = 0, written = 0; written < tot; written += r) {

            if ((r = write(1, buf + written, tot - written)) < 0) {
                perror("write() failed");
                free(buf);
                return;
            }
        }
    }

    end = RDTSC();
    clock_gettime(CLOCK_REALTIME, &ts_after);

    c = (end - start) / iters;
    tot_time_real = timespec_diff(&ts_after, &ts_before);
    tot_time = tot_time_real / iters;
    time_c = tot_time / (double)tot;
    cycles_per_sec = (end - start) / tot_time_real;

    if (tot_time_real <= 0.1) {
        iters *= 10;
        goto retry;
    }

    timespec_to_human_str(tot_time_s, sizeof(tot_time_s), tot_time);
    timespec_to_human_str(c_time_s, sizeof(c_time_s), time_c);
    printf("Term size: %d rows x %d cols\n", w.ws_row, w.ws_col);
    printf("Tot iterations: %d\n\n", iters);
    printf("Screen redraw:         %12llu cycles (%s)\n", c, tot_time_s);
    printf("Avg. character cost: %12llu cycles (%s)\n", c / tot, c_time_s);
    printf("Cycles per sec:        %12.0f cycles/sec\n", cycles_per_sec);
    free(buf);
}
```

# Making libmusl applications to work

Vladislav K. Valtchev (2022)

# Why libmusl?

▶ It made no sense to write a custom libc.

▶ Libmusl produces the smallest binaries (~13 KB for "hello world")

▶ It's actively maintained and widely used in the Embedded Linux world

▶ It's supported by https://toolchains.bootlin.com/

▶ **Uclibc-ng** is more customizable but:

　　▶ Typically produces larger binaries

　　▶ Using a pre-built toolchain means no customization anyway

▶ **Dietlibc** is not well-maintained and has no pre-built toolchains

# Libmusl requires TLS support

- ▶ TLS requires set_thread_area()

Vladislav K. Valtchev (2022)

# Libmusl requires TLS support

▶ TLS requires set_thread_area()

▶ Can we cheat by returning –ENOSYS ? ☺

```
  ┌─crt/crt1.c──────────────────────────────────────────────────────────────────────┐
  │   8            int main();                                                         │
  │   9            weak void _init();                                                  │
  │  10            weak void _fini();                                                  │
  │  11            int __libc_start_main(int (*)(), int, char **,                      │
  │  12                    void (*)(), void(*)(), void(*)());                          │
  │  13                                                                                │
  │  14            void _start_c(long *p)                                              │
  │  15            {                                                                   │
  │  16                    int argc = p[0];                                            │
  │  17                    char **argv = (void *)(p+1);                                │
  │ >18                    __libc_start_main(main, argc, argv, _init, _fini, 0);       │
  │  19            }                                                                   │
  │                                                                                    │
  │                                                                                    │
  │                                                                                    │
  │                                                                                    │
  │                                                                                    │
  │                                                                                    │
  │                                                                                    │
  │                                                                                    │
  │                                                                                    │
  └────────────────────────────────────────────────────────────────────────────────┘
remote Thread 1.1 In: _start_c                              L18    PC: 0x804908f
#0  __init_tp (p=0x804c54c <builtin_tls+12>) at src/env/__init_tls.c:17
#1  0x0804967f in static_init_tls (aux=0xbfffdd8) at src/env/__init_tls.c:149
#2  0x08049270 in __init_libc (envp=0xbffff10, pn=<optimized out>) at src/env/__libc_start_main.c:39
#3  0x0804937d in __libc_start_main (main=0x8049195 <main>, argc=1, argv=0xbffff08) at src/env/__libc_start_main.c:79
#4  0x0804908f in _start_c (p=0xbffff04) at crt/crt1.c:18
#5  0x0804905b in _start ()
(gdb) up
#1  0x0804967f in static_init_tls (aux=0xbfffdd8) at src/env/__init_tls.c:149
#2  0x08049270 in __init_libc (envp=0xbffff10, pn=<optimized out>) at src/env/__libc_start_main.c:39
#3  0x0804937d in __libc_start_main (main=0x8049195 <main>, argc=1, argv=0xbffff08) at src/env/__libc_start_main.c:79
#4  0x0804908f in _start_c (p=0xbffff04) at crt/crt1.c:18
(gdb)
```

```
  src/env/__libc_start_main.c
  69              typedef int lsm2_fn(int (*)(int,char **,char **), int, char **);
  70              static lsm2_fn libc_start_main_stage2;
  71
  72              int __libc_start_main(int (*main)(int,char **,char **), int argc, char **argv)
  73              {
  74                      char **envp = argv+argc+1;
  75
  76                      /* External linkage, and explicit noinline attribute if available,
  77                       * are used to prevent the stack frame used during init from
  78                       * persisting for the entire process lifetime. */
> 79                      __init_libc(envp, argv[0]);
  80
  81                      /* Barrier against hoisting application code or anything using ssp
  82                       * or thread pointer prior to its initialization above. */
  83                      lsm2_fn *stage2 = libc_start_main_stage2;
  84                      __asm__ ( "" : "+r"(stage2) : : "memory" );
  85                      return stage2(main, argc, argv);
  86              }
  87
  88              static int libc_start_main_stage2(int (*main)(int,char **,char **), int argc, char **argv)
  89              {
  90                      char **envp = argv+argc+1;
```
```
remote Thread 1.1 In: __libc_start_main                           L79      PC: 0x804937d
#2  0x08049270 in __init_libc (envp=0xbfffff10, pn=<optimized out>) at src/env/__libc_start_main.c:39
#3  0x0804937d in __libc_start_main (main=0x8049195 <main>, argc=1, argv=0xbfffff08) at src/env/__libc_start_main.c:79
#4  0x0804908f in _start_c (p=0xbfffff04) at crt/crt1.c:18
#5  0x0804905b in _start ()
(gdb) up
#1  0x0804967f in static_init_tls (aux=0xbfffffd8) at src/env/__init_tls.c:149
#2  0x08049270 in __init_libc (envp=0xbfffff10, pn=<optimized out>) at src/env/__libc_start_main.c:39
#3  0x0804937d in __libc_start_main (main=0x8049195 <main>, argc=1, argv=0xbfffff08) at src/env/__libc_start_main.c:79
#4  0x0804908f in _start_c (p=0xbfffff04) at crt/crt1.c:18
(gdb) down
#3  0x0804937d in __libc_start_main (main=0x8049195 <main>, argc=1, argv=0xbfffff08) at src/env/__libc_start_main.c:79
(gdb)
```

```
  src/env/__libc_start_main.c
29                      for (i=0; auxv[i]; i+=2) if (auxv[i]<AUX_CNT) aux[auxv[i]] = auxv[i+1];
30                      __hwcap = aux[AT_HWCAP];
31                      if (aux[AT_SYSINFO]) __sysinfo = aux[AT_SYSINFO];
32                      libc.page_size = aux[AT_PAGESZ];
33
34                      if (!pn) pn = (void*)aux[AT_EXECFN];
35                      if (!pn) pn = "";
36                      __progname = __progname_full = pn;
37                      for (i=0; pn[i]; i++) if (pn[i]=='/') __progname = pn+i+1;
38
>39                     __init_tls(aux);
40                      __init_ssp((void *)aux[AT_RANDOM]);
41
42                      if (aux[AT_UID]==aux[AT_EUID] && aux[AT_GID]==aux[AT_EGID]
43                          && !aux[AT_SECURE]) return;
44
45                      struct pollfd pfd[3] = { {.fd=0}, {.fd=1}, {.fd=2} };
46                      int r =
47          #ifdef SYS_poll
48                      __syscall(SYS_poll, pfd, 3, 0);
49          #else
50                      __syscall(SYS_ppoll, pfd, 3, &(struct timespec){0}, 0, _NSIG/8);
remote Thread 1.1 In: __init_libc                          L39    PC: 0x8049270
(gdb) n
(gdb) s
__copy_tls (mem=0x804c540 <builtin_tls> "") at src/env/__init_tls.c:60
(gdb) s
__init_tp (p=0x804c54c <builtin_tls+12>) at src/env/__init_tls.c:17
(gdb) up
#1  0x0804967f in static_init_tls (aux=0xbffffdd8) at src/env/__init_tls.c:149
#2  0x08049270 in __init_libc (envp=0xbfffff10, pn=<optimized out>) at src/env/__libc_start_main.c:39
#3  0x0804937d in __libc_start_main (main=0x8049195 <main>, argc=1, argv=0xbfffff08) at src/env/__libc_start_main.c:79
(gdb) down
#2  0x08049270 in __init_libc (envp=0xbfffff10, pn=<optimized out>) at src/env/__libc_start_main.c:39
(gdb)
```

```
┌─src/env/__init_tls.c────────────────────────────────────────────────────────────
  139                                     0, libc.tls_size, PROT_READ|PROT_WRITE,
  140                                     MAP_ANONYMOUS|MAP_PRIVATE, -1, 0);
  141                        /* -4095...-1 cast to void * will crash on dereference anyway,
  142                         * so don't bloat the init code checking for error codes and
  143                         * explicitly calling a_crash(). */
  144                    } else {
  145                        mem = builtin_tls;
  146                    }
  147
  148                    /* Failure to initialize thread pointer is always fatal. */
 >149                    if (__init_tp(__copy_tls(mem)) < 0)
  150                        a_crash();
  151            }
  152
  153        weak_alias(static_init_tls, __init_tls);
```
```
remote Thread 1.1 In: static_init_tls                          L149   PC: 0x804967f
__copy_tls (mem=0x804c540 <builtin_tls> "") at src/env/__init_tls.c:60
(gdb) s
__init_tp (p=0x804c54c <builtin_tls+12>) at src/env/__init_tls.c:17
(gdb) up
#1  0x0804967f in static_init_tls (aux=0xbffffdd8) at src/env/__init_tls.c:149
#2  0x08049270 in __init_libc (envp=0xbfffff10, pn=<optimized out>) at src/env/__libc_start_main.c:39
#3  0x0804937d in __libc_start_main (main=0x8049195 <main>, argc=1, argv=0xbfffff08) at src/env/__libc_start_main.c:79
(gdb) down
#2  0x08049270 in __init_libc (envp=0xbfffff10, pn=<optimized out>) at src/env/__libc_start_main.c:39
(gdb) down
#1  0x0804967f in static_init_tls (aux=0xbffffdd8) at src/env/__init_tls.c:149
(gdb)
```

```
 src/env/__init_tls.c
    7              #include "pthread_impl.h"
    8              #include "libc.h"
    9              #include "atomic.h"
   10              #include "syscall.h"
   11
   12         volatile int __thread_list_lock;
   13
   14         int __init_tp(void *p)
   15         {
   16                  pthread_t td = p;
  >17                  td->self = td;
   18                  int r = __set_thread_area(TP_ADJ(p));
   19                  if (r < 0) return -1;
   20                  if (!r) libc.can_do_threads = 1;
   21                  td->detach_state = DT_JOINABLE;
   22                  td->tid = __syscall(SYS_set_tid_address, &__thread_list_lock);
   23                  td->locale = &libc.global_locale;
   24                  td->robust_list.head = &td->robust_list.head;
   25                  td->sysinfo = __sysinfo;
   26                  td->next = td->prev = td;
   27                  return 0;
   28         }
```
```
remote Thread 1.1 In: __init_tp                        L17    PC: 0x8049460
__init_tp (p=0x804c54c <builtin_tls+12>) at src/env/__init_tls.c:17
(gdb) up
#1  0x0804967f in static_init_tls (aux=0xbffffdd8) at src/env/__init_tls.c:149
#2  0x08049270 in __init_libc (envp=0xbfffff10, pn=<optimized out>) at src/env/__libc_start_main.c:39
#3  0x0804937d in __libc_start_main (main=0x8049195 <main>, argc=1, argv=0xbfffff08) at src/env/__libc_start_main.c:79
(gdb) down
#2  0x08049270 in __init_libc (envp=0xbfffff10, pn=<optimized out>) at src/env/__libc_start_main.c:39
(gdb) down
#1  0x0804967f in static_init_tls (aux=0xbffffdd8) at src/env/__init_tls.c:149
(gdb) down
#0  __init_tp (p=0x804c54c <builtin_tls+12>) at src/env/__init_tls.c:17
(gdb)
```

```
┌─src/env/__init_tls.c──────────────────────────────────────────────────────────┐
│    7              #include "pthread_impl.h"                                     │
│    8              #include "libc.h"                                             │
│    9              #include "atomic.h"                                           │
│   10              #include "syscall.h"                                          │
│   11                                                                            │
│   12              volatile int __thread_list_lock;                              │
│   13                                                                            │
│   14              int __init_tp(void *p)                                        │
│   15              {                                                             │
│   16                      pthread_t td = p;                                     │
│   17                      td->self = td;                                        │
│   18                      int r = __set_thread_area(TP_ADJ(p));                 │
│  >19                      if (r < 0) return -1;                                  │
│   20                      if (!r) libc.can_do_threads = 1;                      │
│   21                      td->detach_state = DT_JOINABLE;                       │
│   22                      td->tid = __syscall(SYS_set_tid_address, &__thread_list_lock); │
│   23                      td->locale = &libc.global_locale;                     │
│   24                      td->robust_list.head = &td->robust_list.head;         │
│   25                      td->sysinfo = __sysinfo;                              │
│   26                      td->next = td->prev = td;                             │
│   27                      return 0;                                             │
│   28              }                                                             │
└────────────────────────────────────────────────────────────────────────────────┘
remote Thread 1.1 In: __init_tp                          L19    PC: 0x8049468
#2  0x08049270 in __init_libc (envp=0xbfffff10, pn=<optimized out>) at src/env/__libc_start_main.c:39
#3  0x0804937d in __libc_start_main (main=0x8049195 <main>, argc=1, argv=0xbfffff08) at src/env/__libc_start_main.c:79
(gdb) down
#2  0x08049270 in __init_libc (envp=0xbfffff10, pn=<optimized out>) at src/env/__libc_start_main.c:39
(gdb) down
#1  0x0804967f in static_init_tls (aux=0xbfffffdd8) at src/env/__init_tls.c:149
(gdb) down
#0  __init_tp (p=0x804c54c <builtin_tls+12>) at src/env/__init_tls.c:17
(gdb) n
(gdb) p r
$1 = -38
(gdb)
```

```
 src/env/__init_tls.c
   140                              MAP_ANONYMOUS|MAP_PRIVATE, -1, 0);
   141                     /* -4095...-1 cast to void * will crash on dereference anyway,
   142                      * so don't bloat the init code checking for error codes and
   143                      * explicitly calling a_crash(). */
   144                 } else {
   145                     mem = builtin_tls;
   146                 }
   147
   148                 /* Failure to initialize thread pointer is always fatal. */
   149                 if (__init_tp(__copy_tls(mem)) < 0)
  >150                     a_crash();
   151         }
   152
   153         weak_alias(static_init_tls, __init_tls);
```

```
remote Thread 1.1 In: static_init_tls                        L150  PC: 0x8049686
#1  0x0804967f in static_init_tls (aux=0xbffffdd8) at src/env/__init_tls.c:149
(gdb) down
#0  __init_tp (p=0x804c54c <builtin_tls+12>) at src/env/__init_tls.c:17
(gdb) n
(gdb) p r
$1 = -38
(gdb) s
(gdb) p r
$2 = <optimized out>
(gdb) s
static_init_tls (aux=0xbffffdd8) at src/env/__init_tls.c:150
(gdb)
```

# Sometimes cheating works

Vladislav K. Valtchev (2022)

# Sometimes cheating works

- Sometimes it **doesn't.**

# Sometimes cheating works

- Sometimes it **doesn't.**
- Can we try returning 0 instead and see what happens?

```c
struct user_desc {

    u32 entry_number;
    ulong base_addr;
    u32 limit;

    union {

        struct {
            u32 seg_32bit : 1;          /* Controls GDT_32BIT */
            u32 contents : 2;           /* Controls GDT_ACCESS_DC and GDT_ACCESS_EX */
            u32 read_exec_only : 1;     /* Controls GDT_ACCESS_RW */
            u32 limit_in_pages : 1;     /* Controls GDT_GRAN_4KB */
            u32 seg_not_present : 1;    /* Controls GDT_ACCESS_PRESENT */
            u32 useable : 1;
            u32 ignored : 25;
        };

        u32 flags;
    };
};
```

```
┌─src/env/__init_tls.c─────────────────────────────────────────────────────────────────┐
│     14          int __init_tp(void *p)                                                 │
│     15          {                                                                      │
│     16                  pthread_t td = p;                                              │
│ b+  17                  td->self = td;                                                 │
│  >  18                  int r = __set_thread_area(TP_ADJ(p));                          │
│ b+  19                  if (r < 0) return -1;                                          │
│     20                  if (!r) libc.can_do_threads = 1;                               │
│     21                  td->detach_state = DT_JOINABLE;                                │
│     22                  td->tid = __syscall(SYS_set_tid_address, &__thread_list_lock); │
│     23                  td->locale = &libc.global_locale;                             │
└────────────────────────────────────────────────────────────────────────────────────────┘
┌──────────────────────────────────────────────────────────────────────────────────────┐
│     0x804945d <__init_tp+7>       mov    0x8(%ebp),%ebx                               │
│ b+  0x8049460 <__init_tp+10>      mov    %ebx,(%ebx)                                  │
│     0x8049462 <__init_tp+12>      push   %ebx                                         │
│     0x8049463 <__init_tp+13>      call   0x8049ac4 <__set_thread_area>               │
│ b+> 0x8049468 <__init_tp+18>      add    $0x10,%esp                                   │
│     0x804946b <__init_tp+21>      or     $0xffffffff,%edx                             │
│     0x804946e <__init_tp+24>      test   %eax,%eax                                    │
│     0x8049470 <__init_tp+26>      js     0x80494b5 <__init_tp+95>                     │
│     0x8049472 <__init_tp+28>      jne    0x804947e <__init_tp+40>                     │
│     0x8049474 <__init_tp+30>      movl   $0x1,0x804c640                               │
│     0x804947e <__init_tp+40>      movl   $0x1,0x24(%ebx)                              │
└──────────────────────────────────────────────────────────────────────────────────────┘
remote Thread 1.1 In: __init_tp                          L18    PC: 0x8049468
Breakpoint 8 at 0x8049ae8: file src/thread/i386/__set_thread_area.s, line 19.
(gdb) c
Continuing.

Breakpoint 8, __set_thread_area () at src/thread/i386/__set_thread_area.s:19
(gdb) up
#1  0x08049468 in __init_tp (p=0x804c54c <builtin_tls+12>) at src/env/__init_tls.c:18
(gdb) down
#0  __set_thread_area () at src/thread/i386/__set_thread_area.s:19
(gdb) up
#1  0x08049468 in __init_tp (p=0x804c54c <builtin_tls+12>) at src/env/__init_tls.c:18
(gdb)
```

```
B+  0x8049ae8 <__set_thread_area+36>         test    %eax,%eax
    0x8049aea <__set_thread_area+38>         jne     0x8049aff <__set_thread_area+59>
    0x8049aec <__set_thread_area+40>         mov     (%esp),%edx
    0x8049aef <__set_thread_area+43>         mov     %edx,(%ecx)
   >0x8049af1 <__set_thread_area+45>         lea     0x3(,%edx,8),%edx
    0x8049af8 <__set_thread_area+52>         mov     %edx,%gs
    0x8049afa <__set_thread_area+54>         add     $0x10,%esp
    0x8049afd <__set_thread_area+57>         pop     %ebx
    0x8049afe <__set_thread_area+58>         ret
    0x8049aff <__set_thread_area+59>         mov     %ebx,%ecx
    0x8049b01 <__set_thread_area+61>         xor     %ebx,%ebx
    0x8049b03 <__set_thread_area+63>         xor     %edx,%edx
    0x8049b05 <__set_thread_area+65>         mov     %ebx,(%esp)
    0x8049b08 <__set_thread_area+68>         mov     $0x1,%bl
    0x8049b0a <__set_thread_area+70>         mov     $0x10,%dl
    0x8049b0c <__set_thread_area+72>         mov     $0x7b,%al
    0x8049b0e <__set_thread_area+74>         int     $0x80
    0x8049b10 <__set_thread_area+76>         test    %eax,%eax
    0x8049b12 <__set_thread_area+78>         jne     0x8049afa <__set_thread_area+54>
    0x8049b14 <__set_thread_area+80>         mov     $0x7,%dl
    0x8049b16 <__set_thread_area+82>         inc     %al
    0x8049b18 <__set_thread_area+84>         jmp     0x8049af8 <__set_thread_area+52>
```

In EDX we're supposed to have now the entry number in the GDT. Clearly -1 is invalid.

So now we got an invalid selector now in EDX

```
remote Thread 1.1 In: __set_thread_area                    L23    PC: 0x8049af1
(gdb) p $edx
$1 = -1
(gdb) si
(gdb) p $edx
$2 = -1
(gdb) si
(gdb) p $edx
$3 = -1
(gdb) layout asm
(gdb) p $edx * 8 + 3
$4 = -5
(gdb)
```

```
    0x8049ac4 <__set_thread_area>          push    %ebx
    0x8049ac5 <__set_thread_area+1>        push    $0x51
    0x8049ac7 <__set_thread_area+3>        push    $0xfffff
    0x8049acc <__set_thread_area+8>        pushl   0x10(%esp)
    0x8049ad0 <__set_thread_area+12>       call    0x8049ad5 <__set_thread_area+17>
    0x8049ad5 <__set_thread_area+17>       addl    $0x25fb,(%esp)
    0x8049adc <__set_thread_area+24>       pop     %ecx
    0x8049add <__set_thread_area+25>       mov     (%ecx),%edx
    0x8049adf <__set_thread_area+27>       push    %edx
    0x8049ae0 <__set_thread_area+28>       mov     %esp,%ebx
    0x8049ae2 <__set_thread_area+30>       xor     %eax,%eax
    0x8049ae4 <__set_thread_area+32>       mov     $0xf3,%al
    0x8049ae6 <__set_thread_area+34>       int     $0x80
B+  0x8049ae8 <__set_thread_area+36>       test    %eax,%eax
    0x8049aea <__set_thread_area+38>       jne     0x8049aff <__set_thread_area+59>
    0x8049aec <__set_thread_area+40>       mov     (%esp),%edx
    0x8049aef <__set_thread_area+43>       mov     %edx,(%ecx)
    0x8049af1 <__set_thread_area+45>       lea     0x3(,%edx,8),%edx
   >0x8049af8 <__set_thread_area+52>       mov     %edx,%gs
    0x8049afa <__set_thread_area+54>       add     $0x10,%esp
    0x8049afd <__set_thread_area+57>       pop     %ebx
    0x8049afe <__set_thread_area+58>       ret
    0x8049aff <__set_thread_area+59>       mov     %ebx,%ecx
    0x8049b01 <__set_thread_area+61>       xor     %ebx,%ebx
    0x8049b03 <__set_thread_area+63>       xor     %edx,%edx
    0x8049b05 <__set_thread_area+65>       mov     %ebx,(%esp)
    0x8049b08 <__set_thread_area+68>       mov     $0x1,%bl
    0x8049b0a <__set_thread_area+70>       mov     $0x10,%dl
```

And, of course, here we get a GPF

```
remote Thread 1.1 In: __set_thread_area                          L24    PC: 0x8049af8
(gdb) b *0x8049ae8
Note: breakpoint 4 (disabled) also set at pc 0x8049ae8.
Breakpoint 6 at 0x8049ae8: file src/thread/i386/__set_thread_area.s, line 19.
(gdb) c
Continuing.

Breakpoint 6, __set_thread_area () at src/thread/i386/__set_thread_area.s:19
(gdb) si
(gdb)
```

```
 >0xc0101cd7 <fault13+2>    jmp     0xc0101e02 <asm_fault_entry>
  0xc0101cdc <fault14>      push    $0xe
  0xc0101cde <fault14+2>    jmp     0xc0101e02 <asm_fault_entry>
  0xc0101ce3 <fault15>      push    $0x0
  0xc0101ce5 <fault15+2>    push    $0xf
  0xc0101ce7 <fault15+4>    jmp     0xc0101e02 <asm_fault_entry>
  0xc0101cec <fault16>      push    $0x0
  0xc0101cee <fault16+2>    push    $0x10
  0xc0101cf0 <fault16+4>    jmp     0xc0101e02 <asm_fault_entry>
  0xc0101cf5 <fault17>      push    $0x0
  0xc0101cf7 <fault17+2>    push    $0x11
  0xc0101cf9 <fault17+4>    jmp     0xc0101e02 <asm_fault_entry>
  0xc0101cfe <fault18>      push    $0x0
  0xc0101d00 <fault18+2>    push    $0x12
  0xc0101d02 <fault18+4>    jmp     0xc0101e02 <asm_fault_entry>
  0xc0101d07 <fault19>      push    $0x0
  0xc0101d09 <fault19+2>    push    $0x13
  0xc0101d0b <fault19+4>    jmp     0xc0101e02 <asm_fault_entry>
  0xc0101d10 <fault20>      push    $0x0
  0xc0101d12 <fault20+2>    push    $0x14
  0xc0101d14 <fault20+4>    jmp     0xc0101e02 <asm_fault_entry>
  0xc0101d19 <fault21>      push    $0x0
  0xc0101d1b <fault21+2>    push    $0x15
  0xc0101d1d <fault21+4>    jmp     0xc0101e02 <asm_fault_entry>
  0xc0101d22 <fault22>      push    $0x0
  0xc0101d24 <fault22+2>    push    $0x16
  0xc0101d26 <fault22+4>    jmp     0xc0101e02 <asm_fault_entry>
  0xc0101d2b <fault23>      push    $0x0
```

```
remote Thread 1.1 In: fault13                              L44    PC: 0xc0101cd7
(gdb) file ./build/tilck_unstripped
A program is being debugged already.
Are you sure you want to change the file? (y or n) y
Load new symbol table from "./build/tilck_unstripped"? (y or n) y
Reading symbols from ./build/tilck_unstripped...
Error in re-setting breakpoint 2: No source file named /home/vlad/tilck/toolchain2/i386/musl/src/env/__libc_start_main.c.
Error in re-setting breakpoint 3: No source file named /home/vlad/tilck/toolchain2/i386/musl/src/env/__init_tls.c.
Error in re-setting breakpoint 4: No source file named /home/vlad/tilck/toolchain2/i386/musl/src/thread/i386/__set_thread_area.s.
(gdb)
```

What if we returned 0 and set a valid GDT entry number in **user_desc**, without doing anything else?

```
┌─src/thread/i386/__set_thread_area.s──────────────────────────────────────┐
│   17                mov $243,%al                                          │
│   18                int $128                                              │
│B+ 19                testl %eax,%eax                                       │
│   20                jnz 2f                                                │
│   21                movl (%esp),%edx                                      │
│   22                movl %edx,(%ecx)                                      │
│   23                leal 3(,%edx,8),%edx                                  │
│  >24        3:      movw %dx,%gs                                          │
│   25        1:                                                            │
│   26                addl $16,%esp                                         │
│   27                popl %ebx                                             │
└───────────────────────────────────────────────────────────────────────────┘

┌───────────────────────────────────────────────────────────────────────────┐
│   0x8049ae2 <__set_thread_area+30>      xor    %eax,%eax                   │
│   0x8049ae4 <__set_thread_area+32>      mov    $0xf3,%al                   │
│   0x8049ae6 <__set_thread_area+34>      int    $0x80                       │
│B+ 0x8049ae8 <__set_thread_area+36>      test   %eax,%eax                   │
│   0x8049aea <__set_thread_area+38>      jne    0x8049aff <__set_thread_area+59>
│   0x8049aec <__set_thread_area+40>      mov    (%esp),%edx                 │        ╭──────────────────────────────╮
│   0x8049aef <__set_thread_area+43>      mov    %edx,(%ecx)                 │        │  Now EDX contains a valid GDT │
│   0x8049af1 <__set_thread_area+45>      lea    0x3(,%edx,8),%edx           │        │ selector, 0x23, already used for │
│  >0x8049af8 <__set_thread_area+52>      mov    %edx,%gs  ◄────────────────────────┤         userspace data        │
│   0x8049afa <__set_thread_area+54>      add    $0x10,%esp                  │        ╰──────────────────────────────╯
│   0x8049afd <__set_thread_area+57>      pop    %ebx                        │
│   0x8049afe <__set_thread_area+58>      ret                               │
└───────────────────────────────────────────────────────────────────────────┘
remote Thread 1.1 In: __set_thread_area                    L24    PC: 0x8049af8
(gdb) c
Continuing.

Breakpoint 4, __set_thread_area () at src/thread/i386/__set_thread_area.s:19
(gdb) p $eax
$1 = 0
(gdb) si
(gdb) p $edx
$2 = 4
(gdb) si
(gdb) p/x $edx
$3 = 0x23
(gdb)
```

```
 ┌─src/env/__libc_start_main.c─────────────────────────────────────────────┐
 │  36                      __progname = __progname_full = pn;              │
 │  37                      for (i=0; pn[i]; i++) if (pn[i]=='/') __progname = pn+i+1; │
 │  38                                                                       │
 │  39                      __init_tls(aux);                                 │
 │ >40                      __init_ssp((void *)aux[AT_RANDOM]);  ◄────────── │      We passed __init_tls(aux)!!
 │  41                                                                       │
 │  42                      if (aux[AT_UID]==aux[AT_EUID] && aux[AT_GID]==aux[AT_EGID] │
 │  43                          && !aux[AT_SECURE]) return;                  │
 │  44                                                                       │
 │  45                      struct pollfd pfd[3] = { {.fd=0}, {.fd=1}, {.fd=2} }; │
 │  46                      int r =                                          │
 └──────────────────────────────────────────────────────────────────────────┘
 ┌──────────────────────────────────────────────────────────────────────────┐
 │    0x804925f <__init_libc+154>      jmp     0x804924d <__init_libc+136>   │
 │    0x8049261 <__init_libc+156>      sub     $0xc,%esp                     │
 │    0x8049264 <__init_libc+159>      lea     -0xb0(%ebp),%eax              │
 │    0x804926a <__init_libc+165>      push    %eax                          │
 │    0x804926b <__init_libc+166>      call    0x8049529 <static_init_tls>   │
 │   >0x8049270 <__init_libc+171>      pop     %eax                          │
 │    0x8049271 <__init_libc+172>      pushl   -0x4c(%ebp)                   │
 │    0x8049274 <__init_libc+175>      call    0x80491c4 <dummy1>            │
 │    0x8049279 <__init_libc+180>      add     $0x10,%esp                    │
 │    0x804927c <__init_libc+183>      mov     -0x80(%ebp),%eax              │
 │    0x804927f <__init_libc+186>      cmp     %eax,-0x84(%ebp)              │
 │    0x8049285 <__init_libc+192>      jne     0x8049295 <__init_libc+208>   │
 └──────────────────────────────────────────────────────────────────────────┘
remote Thread 1.1 In: __init_libc                          L40    PC: 0x8049270
(gdb) c
Continuing.

Breakpoint 5, 0x08049493 in __syscall1 (a1=134530688, n=258) at ./arch/i386/syscall_arch.h:25
(gdb) si
__init_tp (p=0x804c54c <builtin_tls+12>) at ./arch/i386/syscall_arch.h:26
0x0804967f in static_init_tls (aux=0xbffffdd8) at src/env/__init_tls.c:149
(gdb) si
(gdb) p $eax
$4 = 0
(gdb) si
__init_libc (envp=0xbfffff10, pn=<optimized out>) at src/env/__libc_start_main.c:40
(gdb)
```

```
  src/env/__libc_start_main.c
   86        }
   87
   88        static int libc_start_main_stage2(int (*main)(int,char **,char **), int argc, char **argv)
   89        {
  >90                char **envp = argv+argc+1;
   91                __libc_start_init();
   92
   93                /* Pass control to the application */
   94                exit(main(argc, argv, envp));
   95                return 0;
   96        }
```

```
>0x8049338 <libc_start_main_stage2>       push    %ebp
 0x8049339 <libc_start_main_stage2+1>     mov     %esp,%ebp
 0x804933b <libc_start_main_stage2+3>     push    %edi
 0x804933c <libc_start_main_stage2+4>     push    %esi
 0x804933d <libc_start_main_stage2+5>     push    %ebx
 0x804933e <libc_start_main_stage2+6>     sub     $0xc,%esp
 0x8049341 <libc_start_main_stage2+9>     mov     0xc(%ebp),%ebx
 0x8049344 <libc_start_main_stage2+12>    mov     0x10(%ebp),%esi
 0x8049347 <libc_start_main_stage2+15>    call    0x8049316 <libc_start_init>
 0x804934c <libc_start_main_stage2+20>    lea     0x4(%esi,%ebx,4),%edi
 0x8049350 <libc_start_main_stage2+24>    push    %eax
 0x8049351 <libc_start_main_stage2+25>    push    %edi
```

```
remote Thread 1.1 In: libc_start_main_stage2                    L90    PC: 0x8049338
0x0804967f in static_init_tls (aux=0xbffffdd8) at src/env/__init_tls.c:149
(gdb) si
(gdb) p $eax
$4 = 0
(gdb) si
__init_libc (envp=0xbfffff10, pn=<optimized out>) at src/env/__libc_start_main.c:40
(gdb) si
dummy1 (p=0x0) at src/env/__libc_start_main.c:15
__init_libc (envp=0xbfffff10, pn=<optimized out>) at src/env/__libc_start_main.c:42
__libc_start_main (main=0x8049195 <main>, argc=1, argv=0xbfffff08) at src/env/__libc_start_main.c:84
(gdb) si
libc_start_main_stage2 (main=0x8049195 <main>, argc=1, argv=0xbfffff08) at src/env/__libc_start_main.c:90
(gdb)
```

```
./arch/i386/syscall_arch.h
    36          static inline long __syscall3(long n, long a1, long a2, long a3)
    37          {
    38                  unsigned long __ret;
    39          #if !defined(__PIC__) || !defined(BROKEN_EBX_ASM)
  > 40                  __asm__ __volatile__ (SYSCALL_INSNS : "=a"(__ret) : "a"(n), "b"(a1), "c"(a2), "d"(a3) : "memory");
    41          #else
    42                  __asm__ __volatile__ (SYSCALL_INSNS_34 : "=a"(__ret) : "a"(n), "D"(a1), "c"(a2), "d"(a3) : "memory");
    43          #endif
    44                  return __ret;
    45          }
    46


    0x8049840 <__stdout_write+8>      mov     0x8(%ebp),%esi
    0x8049843 <__stdout_write+11>     movl    $0x8049c1d,0x24(%esi)
    0x804984a <__stdout_write+18>     testb   $0x40,(%esi)
    0x804984d <__stdout_write+21>     jne     0x8049871 <__stdout_write+57>
    0x804984f <__stdout_write+23>     lea     -0x10(%ebp),%edx
    0x8049852 <__stdout_write+26>     mov     $0x36,%eax
    0x8049857 <__stdout_write+31>     mov     $0x5413,%ecx
    0x804985c <__stdout_write+36>     mov     0x3c(%esi),%ebx
  > 0x804985f <__stdout_write+39>     call    *%gs:0x10
    0x8049866 <__stdout_write+46>     test    %eax,%eax
    0x8049868 <__stdout_write+48>     je      0x8049871 <__stdout_write+57>
    0x804986a <__stdout_write+50>     movl    $0xffffffff,0x50(%esi)
```

> Ehm.. I don't believe we're going to pass that far indirect call…

```
remote Thread 1.1 In: __stdout_write                               L40    PC: 0x804985f
#2  0x080497de in __overflow (f=0x804c040 <__stdout_FILE>, _c=10) at src/stdio/__overflow.c:8
#3  0x0804942d in puts (s=0x804a000 "hello world") at src/stdio/puts.c:7
#4  0x080491b3 in main () at userapps/hello.c:5
(gdb) down
#0  0x0804985c in __syscall3 (a3=-1073742344, a2=21523, a1=1, n=54) at ./arch/i386/syscall_arch.h:40
(gdb) bt
#0  0x0804985c in __syscall3 (a3=-1073742344, a2=21523, a1=1, n=54) at ./arch/i386/syscall_arch.h:40
#1  __stdout_write (f=0x804c040 <__stdout_FILE>, buf=0xbffffe2f "\n", len=1) at src/stdio/__stdout_write.c:8
#2  0x080497de in __overflow (f=0x804c040 <__stdout_FILE>, _c=10) at src/stdio/__overflow.c:8
#3  0x0804942d in puts (s=0x804a000 "hello world") at src/stdio/puts.c:7
#4  0x080491b3 in main () at userapps/hello.c:5
(gdb) si
(gdb)
```

```
┌─/home/vlad/tilck/kernel/arch/i386/fault_handlers.S──────────────────────────────┐
│  >45           fault_with_err_code 14 # Page Fault Exception                      │
│   46                                                                              │
│   47           fault 15 # Reserved Exception                                      │
│   48           fault 16 # Floating Point Exception                                │
│   49           fault 17 # Alignment Check Exception                               │
│   50           fault 18 # Machine Check Exception                                 │
│   51                                                                              │
│   52           fault 19                                                           │
│   53           fault 20                                                           │
│   54           fault 21                                                           │
│   55           fault 22                                                           │
│                                                                                   │
└───────────────────────────────────────────────────────────────────────────────────┘
```

```
 >0xc0101cde <fault14+2>  jmp     0xc0101e02 <asm_fault_entry>
  0xc0101ce3 <fault15>     push    $0x0
  0xc0101ce5 <fault15+2>   push    $0xf
  0xc0101ce7 <fault15+4>   jmp     0xc0101e02 <asm_fault_entry>
  0xc0101cec <fault16>     push    $0x0
  0xc0101cee <fault16+2>   push    $0x10
  0xc0101cf0 <fault16+4>   jmp     0xc0101e02 <asm_fault_entry>
  0xc0101cf5 <fault17>     push    $0x0
  0xc0101cf7 <fault17+2>   push    $0x11
  0xc0101cf9 <fault17+4>   jmp     0xc0101e02 <asm_fault_entry>
  0xc0101cfe <fault18>     push    $0x0
  0xc0101d00 <fault18+2>   push    $0x12
```

Yep, page fault.

```
remote Thread 1.1 In: fault14                        L45    PC: 0xc0101cde
0xc0101cde in ?? ()
(gdb) p $eax
$6 = 54
(gdb) p/x $eax
$7 = 0x36
(gdb) file ./build/tilck_unstripped
A program is being debugged already.
Are you sure you want to change the file? (y or n) y
Load new symbol table from "./build/tilck_unstripped"? (y or n) y
Reading symbols from ./build/tilck_unstripped...
Error in re-setting breakpoint 2: No source file named __init_tls.c.
Error in re-setting breakpoint 3: No source file named __set_thread_area.s.
(gdb) |
```

```
/home/vlad/tilck/kernel/arch/i386/paging.c
224             panic("PAGE FAULT in attempt to %s %p from %s%s\nEIP: %p [%s + %d]\n",
225                 rw ? "WRITE" : "READ",
226                 vaddr,
227                 "kernel",
228                 !p ? " (NON present)." : ".",
229                 r->eip, sym_name ? sym_name : "???", off);
230         }
231
232     void handle_page_fault_int(regs_t *r)
233     {
234         u32 vaddr;
235         asmVolatile("movl %%cr2, %0" : "=r"(vaddr));
236
237         bool p  = !!(r->err_code & PAGE_FAULT_FL_PRESENT);
238         bool rw = !!(r->err_code & PAGE_FAULT_FL_RW);
239         bool us = !!(r->err_code & PAGE_FAULT_FL_US);
>240         int sig = SIGSEGV;
241         struct user_mapping *um;
242
243         if (!us) {
244             /*
245              * Tilck does not support kernel-space page faults caused by the kernel,
246              * while it allows user-space page faults caused by kernel (CoW pages).
247              * Therefore, such a fault is necessary caused by a bug.
```

remote Thread 1.1 In: handle_page_fault_int                          L240   PC: 0xc0103e55
(gdb) n
(gdb) s
handle_page_fault_int (r=0xf8032fa8) at /home/vlad/tilck/kernel/arch/i386/paging.c:235
(gdb) n
(gdb) p p
$8 = false
(gdb) p rw
$9 = false
(gdb) p us
$10 = true
(gdb) p/x vaddr
$11 = 0x10
(gdb)

Vaddr is clearly just 0x10 because the GDT
selector 0x23 has offset = 0 (flat segmentation)

# Lesson learned

- Often, we cannot cheat.

# Lesson learned

- Often, we cannot cheat.
- Even basic I/O functions use TLS variables.

# Lesson learned

▶ Often, we cannot cheat.

▶ Even basic I/O functions use TLS variables.

▶ Had to provide a fully-functional implementation for **set_thread_area()**, in order run even single-threaded libmusl applications.

```c
384  int sys_set_thread_area(void *arg)
385  {
386      int rc = 0;
387      struct gdt_entry e = {0};
388      struct user_desc dc;
389      struct user_desc *ud = arg;
390
391      rc = copy_from_user(&dc, ud, sizeof(struct user_desc));
392      if (rc != 0)
393          return -EFAULT;
394
395
396      disable_preemption();
397
398      if (!(dc.flags == USER_DESC_FLAGS_EMPTY && !dc.base_addr && !dc.limit)) {
399          gdt_set_entry(&e, dc.base_addr, dc.limit, 0, 0);
400          e.s = 1;
401          e.dpl = 3;
402          e.d = dc.seg_32bit;
403          e.type |= (dc.contents << 2);
404          e.type |= !dc.read_exec_only ? GDT_ACCESS_RW : 0;
405          e.g = dc.limit_in_pages;
406          e.avl = dc.useable;
407          e.p = !dc.seg_not_present;
408      } else {
409          /* The user passed an empty descriptor: entry_number cannot be -1 */
410          if (dc.entry_number == INVALID_ENTRY_NUM) {
411              rc = -EINVAL;
412              goto out;
413          }
414      }
415
416      if (dc.entry_number == INVALID_ENTRY_NUM) {
417
418          int slot = find_available_slot_in_user_task();
419          if (slot < 0) {
420              rc = -ESRCH;
421              goto out;
422          }
423
424          dc.entry_number = (u32)gdt_add_entry(&e);
425          if (dc.entry_number == INVALID_ENTRY_NUM) {
426
427              rc = gdt_expand();
428
429              if (rc < 0) {
430                  rc = -ESRCH;
431                  goto out;
432              }
433              dc.entry_number = (u32)gdt_add_entry(&e);
434              ASSERT(dc.entry_number != INVALID_ENTRY_NUM);
435          }
436          gdt_set_slot(get_curr_proc(), (u16)slot, (u16)dc.entry_number);
437          goto out;
438      }
439
440      /* Handling the case where the user specified a GDT entry number */
441
442      int slot = get_user_task_slot_for_gdt_entry(dc.entry_number);
443
444      if (slot < 0) {
445          /* A GDT entry with that index has never been allocated by this task */
446
447          if (dc.entry_number >= gdt_size || gdt[dc.entry_number].access) {
448              /* The entry is out-of-bounds or it's used by another task */
449              rc = -EINVAL;
450              goto out;
451          }
452
453          /* The entry is available, now find a slot */
454          slot = find_available_slot_in_user_task();
455
456          if (slot < 0) {
457              /* Unable to find a free slot in this struct task struct */
458              rc = -ESRCH;
459              goto out;
460          }
461
462          gdt_set_slot(get_curr_proc(), (u16)slot, (u16)dc.entry_number);
463      }
464
465      ASSERT(dc.entry_number < gdt_size);
466
467      set_entry_num(dc.entry_number, &e);
468
469      /*
470       * We're here because either we found a slot already containing this index
471       * (therefore it must be valid) or the index is in-bounds and it is free.
472       */
473
474  out:
475      enable_preemption();
476
477      if (!rc) {
478
479          /*
480           * Positive case: we get here with rc = SUCCESS, now flush back the
481           * the struct user_desc (we might have changed its entry_number).
482           */
483          rc = copy_to_user(ud, &dc, sizeof(struct user_desc));
484
485          if (rc < 0)
486              rc = -EFAULT;
487      }
488
489      return rc;
490  }
```

That was quite some code, but it's not enough. We need a **ref-count** for GDT entries as well.

Why? Think about **fork().** What happens if the parent dies before the child and we free the GDT slots?

```
780 void
781 arch_specific_new_proc_setup(struct process *pi, struct process *parent)
782 {
783     arch_proc_members_t *arch = get_proc_arch_fields(pi);
784
785     if (!parent)
786         return;        /* we're done */
787
788     memcpy(&pi->pi_arch, &parent->pi_arch, sizeof(pi->pi_arch));
789
790     if (arch->ldt)
791         gdt_entry_inc_ref_count(arch->ldt_index_in_gdt);
792
793     for (int i = 0; i < ARRAY_SIZE(arch->gdt_entries); i++)
794         if (arch->gdt_entries[i])
795             gdt_entry_inc_ref_count(arch->gdt_entries[i]);
796
797     pi->set_child_tid = NULL;
798 }
```

```
800 void
801 arch_specific_free_proc(struct process *pi)
802 {
803     arch_proc_members_t *arch = get_proc_arch_fields(pi);
804
805     if (arch->ldt) {
806         gdt_clear_entry(arch->ldt_index_in_gdt);
807         arch->ldt = NULL;
808     }
809
810     for (int i = 0; i < ARRAY_SIZE(arch->gdt_entries); i++) {
811         if (arch->gdt_entries[i]) {
812             gdt_clear_entry(arch->gdt_entries[i]);
813             arch->gdt_entries[i] = 0;
814         }
815     }
816 }
```

# ACPICA & AcpiOsWaitSemaphore()

- ACPICA requires the OSL to provide a counting semaphore implementation capable of waiting and signaling N units.

- That is weird requirement.

- It could be trivially implemented on the top of a regular counting semaphore, but that would be extremely inefficient.

- I implemented such a semaphore in Tilck.

# Classic semaphore

```c
20  void ksem_wait(struct ksem *s)
21  {
22      struct task *curr = get_curr_task();
23      disable_preemption();
24
25      if (--s->counter < 0) {
26
27          task_set_wait_obj(curr, WOBJ_SEM, s, NO_EXTRA, &s->wait_list);
28          enable_preemption_nosched();
29          kernel_yield();
30          return;
31      }
32
33      enable_preemption();
34  }
```

# New semaphore [1/2]

```c
25  int ksem_wait(struct ksem *s, int units, int timeout_ticks)
26  {
27      int rc = -ETIME;
28      ASSERT(units > 0);
29
30      if (s->max != KSEM_NO_MAX && units > s->max)
31          return -EINVAL;
32
33      disable_preemption();
34      {
35          if (timeout_ticks != KSEM_NO_WAIT) {
36              u64 start_ticks, end_ticks;
37
38              if (timeout_ticks > 0) {
39
40                  start_ticks = get_ticks();
41                  end_ticks = start_ticks + (u32)timeout_ticks;
42
43                  if (s->counter < units)
44                      task_set_wakeup_timer(get_curr_task(), (u32)timeout_ticks);
45              }
46
47              while (s->counter < units) {
48
49                  if (timeout_ticks > 0 && get_ticks() >= end_ticks)
50                      break;
51
52                  prepare_to_wait_on(WOBJ_SEM, s, (u32)units, &s->wait_list);
53                  enter_sleep_wait_state(); /* after that, preemption will be enabled */
54                  disable_preemption();
55              }
56
57              if (timeout_ticks > 0)
58                  task_cancel_wakeup_timer(get_curr_task());
59          }
60
61          if (s->counter >= units) {
62              s->counter -= units;
63              rc = 0;
64          }
65      }
66      enable_preemption();
67      return rc;
68  }
```

# Classic semaphore

```
37 void ksem_signal(struct ksem *s)
38 {
39     disable_preemption();
40
41     if (s->counter++ < 0) {
42
43         ASSERT(!list_is_empty(&s->wait_list));
44
45         struct wait_obj *task_wo =
46             list_first_obj(&s->wait_list, struct wait_obj, wait_list_node);
47
48         struct task *ti = CONTAINER_OF(task_wo, struct task, wobj);
49         task_reset_wait_obj(ti);
50     }
51
52     enable_preemption();
53 }
```

# New semaphore [2/2]

```
70 int ksem_signal(struct ksem *s, int units)
71 {
72     struct wait_obj *wo, *tmp;
73     int rem_counter, rc = 0;
74     ASSERT(units > 0);
75
76     disable_preemption();
77
78     if (s->max != KSEM_NO_MAX) {
79
80         if (units > s->max) {
81             rc = -EINVAL;
82             goto out;
83         }
84
85         if (s->counter > s->max - units) {
86             rc = -EDQUOT;
87             goto out;
88         }
89     }
90
91     s->counter += units;
92     rem_counter = s->counter;
93
94     list_for_each(wo, tmp, &s->wait_list, wait_list_node) {
95
96         if (rem_counter <= 0)
97             break; /* not enough units to unblock anybody */
98
99         int wait_units = (int)wo->extra;
100
101         if (wait_units <= rem_counter) {
102             struct task *ti = CONTAINER_OF(wo, struct task, wobj);
103             rem_counter -= wait_units;
104             wake_up(ti);
105         }
106     }
107
108 out:
109     enable_preemption();
110     return rc;
111 }
```

# But.. how Linux did implement the counting semaphore to make ACPICA happy?

Vladislav K. Valtchev (2022)

# But.. how Linux did implement the counting semaphore to make ACPICA happy?

It didn't ☺

# But.. how Linux did implement the counting semaphore to make ACPICA happy?

```
1239    /*
1240     * TODO: Support for units > 1?
1241     */
1242    acpi_status acpi_os_wait_semaphore(acpi_handle handle, u32 units, u16 timeout)
1243    {
1244            acpi_status status = AE_OK;
1245            struct semaphore *sem = (struct semaphore *)handle;
1246            long jiffies;
1247            int ret = 0;
1248
1249            if (!acpi_os_initialized)
1250                    return AE_OK;
1251
1252            if (!sem || (units < 1))
1253                    return AE_BAD_PARAMETER;
1254
1255            if (units > 1)
1256                    return AE_SUPPORT;
1257
1258            ACPI_DEBUG_PRINT((ACPI_DB_MUTEX, "Waiting for semaphore[%p|%d|%d]\n",
1259                            handle, units, timeout));
```

# But.. how Linux did implement the counting semaphore to make ACPICA happy?

```
1239     /*
1240      * TODO: Support for units > 1?
1241      */
1242    acpi_status acpi_os_wait_semaphore(acpi_handle handle, u32 units, u16 timeout)
1243    {
1244            acpi_status status = AE_OK;
1245            struct semaphore *sem = (struct semaphore *)handle;
1246            long jiffies;
1247            int ret = 0;
1248
1249            if (!acpi_os_initialized)
1250                    return AE_OK;
1251
1252            if (!sem || (units < 1))
1253                    return AE_BAD_PARAMETER;
1254
1255            if (units > 1)
1256                    return AE_SUPPORT;
1257
1258            ACPI_DEBUG_PRINT((ACPI_DB_MUTEX, "Waiting for semaphore[%p|%d|%d]\n",
1259                             handle, units, timeout));
```

*Sometimes cheating works.*