

Kernel Rookie Guide Workshop

Andrzej Pietrasiewicz
andrzej.p@samsung.com

Samsung R&D Institute Poland
Warsaw

Open Source Summit Europe
21st October 2018

SAMSUNG



Kernel Rookie
Guide Workshop

Andrzej
Pietrasiewicz

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Requirements 1/3

- ▶ Some C knowledge
- ▶ Ability to use 'make'
- ▶ Own laptop with Linux available
- ▶ Network connectivity
- ▶ Internet browser

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Requirements 2/3

▶ In Linux:

- ▶ qemu-system-arm version at least 2.5.0
- ▶ Cross toolchain ¹
- ▶ Tools for building the kernel
- ▶ tar
- ▶ wget
- ▶ editor
- ▶ `wget https://mirrors.edge.kernel.org/pub/linux/kernel/v4.x/linux-4.18.tar.xz2`
- ▶ `workshops.tar.gz` (from the instructor)
- ▶ extract `linux-4.18.tar.xz` and `workshops tar.gz` next to each other (expect `linux-4.18` and `workshops` directories next to each other)

¹arm-linux-gnueabi-gcc expected in \$PATH

²`tar xxf <filename>`

Requirements 3/3

- ▶ On a Ubuntu 18.04 Desktop fresh install³
 - ▶ `sudo apt-get install qemu-system-arm make gcc bison flex pkg-config libncurses-dev gcc-arm-linux-gnueabi vim-gtk`
- ▶ On older Ubuntu you might also need
 - ▶ `sudo apt-get install bc libssl-dev`
- ▶ Translate these requirements to the distro you intend to use at workshops!

³I tried it in a virtual machine

Workshop plan

Meeting the Linux kernel

- OS purpose and Linux Hardware, Kernel & toolchain (Early) userspace

Writing basic modules

- In-tree vs out-of-tree
- Module boilerplate code
- Linked lists
- Small allocations

Userspace interface

- Syscalls
- Filesystems
- Devices

Synchronization and locking basics

- Concurrency
- Synchronization primitives
- Final workshop

What now?

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls

Filesystems

Devices

Synchronization and locking basics

Concurrency

Synchronization primitives

Final workshop

What now?

Kernel Rookie
Guide Workshop

Andrzej
Pietrasiewicz

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Workshop plan

Meeting the Linux kernel

- OS purpose and Linux

- Hardware, Kernel & toolchain

- (Early) userspace

Writing basic modules

- In-tree vs out-of-tree

- Module boilerplate code

- Linked lists

- Small allocations

Userspace interface

- Syscalls

- Filesystems

- Devices

Synchronization and locking basics

- Concurrency

- Synchronization primitives

- Final workshop

What now?

Rationale

- ▶ Ancient/bare-metal programming vs. modern programming
- ▶ More tasks than CPUs at a time (some vs. 1 or more)
- ▶ Scarce resources (cannot satisfy everyone at the same time)
- ▶ Devices

OS Purpose

- ▶ Manage processes (give everyone a chance to execute)
- ▶ Manage memory (code and data require memory)
- ▶ Manage devices (access peripherals)

OS major design decisions

- ▶ Microkernel vs monolithic (Tannenbaum vs Torvalds)
- ▶ How to use hardware capabilities
- ▶ How to interoperate with users
- ▶ Development model & project governance

- ▶ Monolithic, but with loadable modules
- ▶ Development discussed on mailing lists, patches
- ▶ Source code in git^{4,5}
- ▶ ~800MB w/o version control
- ▶ Indexing, also on-line, e.g.: <https://elixir.bootlin.com>

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

⁴[git://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git](https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git)

⁵<https://git-scm.com/>

Workshop 0⁶

- ▶ Find *container_of* definition in 4.18 and try understanding it
- ▶ If more than one found choose the one whose location looks most general
- ▶ <https://elixir.bootlin.com>, use 4.18

⁶Look for `TODO.txt` in `W0` from the extracted `tar.gz`

Workshop 0 solution

Kernel Rookie
Guide Workshop

Andrzej
Pietrasiewicz

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization primitives
Final workshop

What now?

- ▶ Originally i386
- ▶ Many added, incl. arm and arm64
- ▶ Hardware discovery (ACPI, dtb⁸, autodetection)
- ▶ Workshops
 - ▶ Emulated hardware - qemu⁹
 - ▶ ARM Versatile Express (qemu: vexpress-a9)

⁸<https://www.devicetree.org/>

⁹<https://www.qemu.org/>

Kernel configuration

- ▶ Highly configurable
 - ▶ Architectures
 - ▶ Kernel subsystems (/kernel, /mm, /net, /drivers/*, /arch/* etc.)
 - ▶ Peripherals
 - ▶ Compile-in/Modularize/Leave out
- ▶ Dedicated tools complementing the Makefile

Kconfig language for specifying kernel composition and dependencies

defconfig sane set of choices for given architecture

.config result of defconfig application

```
1 # sane set of settings
2 make ARCH=arm vexpress_defconfig
3
4 # customization if needed
5 make ARCH=arm menuconfig
```

Listing 1: Configuring the kernel

Kernel building

- ▶ (Cross-)toolchain¹⁰ (ARCH=..., CROSS_COMPILE=...)
- ▶ Keep your CPU cores busy with -j!

```
1 # either
2 make -j8 ARCH=arm CROSS_COMPILE=arm-linux-gnueabi-
3
4 # or
5 export CROSS_COMPILE=arm-linux-gnueabi-
6 make -j8 ARCH=arm
```

Listing 2: Compiling the kernel

- ▶ Make targets

```
1 make ARCH=arm help
```

Listing 3: Getting help on make targets

¹⁰arm-linux-gnueabi-gcc expected in \$PATH

Workshop 1¹¹

- ▶ Patch the kernel with the provided patches
- ▶ Configure & cross-compile the kernel for vexpress
- ▶ Analyze the vexpress-run.sh script
- ▶ Run emulated machine in qemu with the script
 - ▶ ^A-X to exit qemu

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

¹¹Look for TODO.txt in W1

Workshop 1 solution

Kernel Rookie
Guide Workshop

Andrzej
Pietrasiewicz

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization primitives
Final workshop

What now?

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

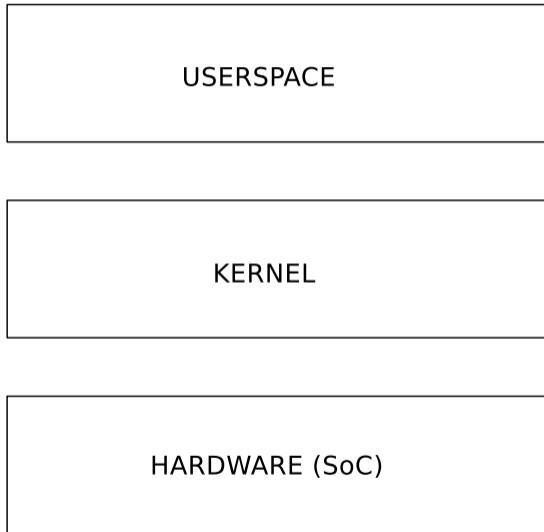
Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

System layers



Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Ramfs, tmpfs, rootfs, initramfs

ramfs disk caches mounted as filesystem, no backing storage

tmpfs ramfs += size limit, swapping

rootfs special instance of ramfs/tmpfs, always present

initramfs optional cpio archive¹², extracted into rootfs at boot time

- ▶ Initramfs handy for experiments with kernel
- ▶ We use external initramfs for workshops

¹²built into the kernel and/or provided by bootloader

Minimal userspace for workshops

- ▶ busybox¹³ + minimal /etc packaged into initramfs
- ▶ Single binary containing 'applets' corresponding to standard tools
- ▶ Tool from kernel source tree understands both directories and list files

```
1 cat initramfs.list
2 file /busybox-armv7r ../initramfs/busybox-armv7r 0755 0 0
3 dir /sbin 0755 0 0
4 dir /proc 0755 0 0
5 dir /sys 0755 0 0
6 dir /bin 0755 0 0
7 dir /usr 0755 0 0
8 dir /usr/sbin 0755 0 0
9 dir /usr/bin 0755 0 0
10 dir /etc 0755 0 0
11 dir /etc/init.d 0755 0 0
12 file /etc/inittab ../initramfs/inittab 0755 0 0
13 file /etc/init.d/rcS ../initramfs/rcS 0755 0 0
```

Listing 4: Initramfs list - host paths relative to kernel source topdir!

¹³<https://busybox.net/>

Workshop 2¹⁴

- ▶ Analyze `generate_initramfs.sh`
- ▶ Add a short text file to `initramfs` image using `generate_initramfs.sh`
- ▶ Run the virtual machine and verify new file's presence and contents

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Workshop 2 solution

Kernel Rookie
Guide Workshop

Andrzej
Pietrasiewicz

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization primitives
Final workshop

What now?

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization primitives
Final workshop

What now?

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Source code placement

in-tree inside kernel source directory

out-of-tree outside kernel source directory

- ▶ Kernel build system happily handles both!

Meeting the Linux kernel
OS purpose and Linux Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules
In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface
Syscalls
Filesystems
Devices

Synchronization and locking basics
Concurrency
Synchronization primitives
Final workshop

What now?

Out-of-tree Makefile

- ▶ LDD3¹⁵ provides a nice Makefile (meant for compiling modules for currently running kernel)
- ▶ The below code creates hello.ko module (from hello.o and implicitly from hello.c)

```
1 # If KERNELRELEASE is defined, we've been invoked from the
2 # kernel build system and can use its language.
3 ifneq ($(KERNELRELEASE),)
4     obj-m := hello.o
5 # Otherwise we were called directly from the command
6 # line; invoke the kernel build system.
7 else
8     KERNELDIR ?= /lib/modules/$(shell uname -r)/build
9     PWD := $(shell pwd)
10 default:
11     $(MAKE) -C $(KERNELDIR) M=$(PWD) modules
12 endif
```

Listing 5: Out-of-tree module makefile

Interaction with modules

- ▶ Module dependencies
- ▶ Userspace tools
 - `insmod` dependencies must be provided
 - `modprobe` dependencies resolved
 - `rmmod` unload only this module
 - `modprobe -r` unload this module and unused dependencies
 - `lsmod` list loaded modules
 - `depmod` analyze module dependencies

```
1 # loading
2 insmod <filename>.ko
3 # or
4 modprobe <module name>
5
6 #unloading
7 rmmod <module name>
8 # or
9 modprobe -r <module name>
10
11 # list loaded
12 lsmod
```

Listing 6: Loading/unloading kernel modules

Workshop plan

Meeting the Linux kernel

- OS purpose and Linux Hardware, Kernel & toolchain (Early) userspace

Writing basic modules

- In-tree vs out-of-tree
- Module boilerplate code**
- Linked lists
- Small allocations

Userspace interface

- Syscalls
- Filesystems
- Devices

Synchronization and locking basics

- Concurrency
- Synchronization primitives
- Final workshop

What now?

Module boilerplate code

```
1 #include <linux/init.h>
2 #include <linux/module.h>
3
4 static int hello_init(void)
5 {
6     return 0;
7 }
8 module_init(hello_init);
9
10 static void hello_exit(void)
11 {
12 }
13 module_exit(hello_exit);
14
15 MODULE_LICENSE("GPL");
```

Listing 7: Linux kernel module boilerplate code

- ▶ A value passed to module at insertion^{16,17}
- ▶ Several types (int, string, etc)
- ▶ A way of associating the said value with a variable in the module

Meeting the Linux kernel

OS purpose and Linux Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization primitives
Final workshop

What now?

```
1 /* use this at file scope! */  
2 module_param(<var_name>, int, <flags>) /* flags for sysfs access mode */
```

Listing 8: Module parameters (selection)

¹⁶or in the kernel command line if module compiled-in

¹⁷or through sysfs if parameter available there

Kernel diagnostic messages circular buffer

- ▶ Always available
- ▶ Never runs out of space¹⁸
- ▶ Some messages can be simultaneously output to the current system console¹⁹
- ▶ `dmesg` to output the diagnostic messages circular buffer contents

¹⁸But older messages can be overwritten

¹⁹Depends on message levels and current kernel settings

Message levels

- ▶ message levels KERN_*

EMERG (0) system is unusable

ALERT (1) action must be taken immediately

CRIT (2) critical conditions

ERR (3) error conditions

WARNING (4) warning conditions

NOTICE (5) normal but significant condition

INFO (6) informational

DEBUG (7) debug-level messages

```
1 cat /proc/sys/kernel/printk  
2 4 4 1 7
```

Listing 9: Printk levels

Writing messages to the console

▶ Output methods

`printk` `printf` equivalent

`pr_<level>` `printf` equivalent (level: emerg, alert, crit etc.)²⁰

▶ Remember about newlines!

```
1 /* default level */
2 printk("Hello, \u00a0world!\n");
3
4 /* selected level */
5 printk(KERN_WARNING "Hello, \u00a0warning!\n");
6
7 pr_warning("Hello, \u00a0warning!\n");
```

Listing 10: Outputting messages from kernel

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code

Linked lists

Small allocations

Userspace interface

Syscalls

Filesystems

Devices

Synchronization
and locking basics

Concurrency

Synchronization
primitives

Final workshop

What now?

- ▶ Write, compile and deploy (initramfs!) a module which does (almost) nothing
 - ▶ Module's init shall output a welcome string
 - ▶ module's exit shall output a goodbye string
- ▶ Load/unload the module in the virtual machine

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code

Linked lists

Small allocations

Userspace interface

Syscalls

Filesystems

Devices

Synchronization
and locking basics

Concurrency

Synchronization
primitives

Final workshop

What now?

Workshop 3 solution

Kernel Rookie
Guide Workshop

Andrzej
Pietrasiewicz

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
**Module boilerplate
code**
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code

Linked lists

Small allocations

Userspace interface

Syscalls

Filesystems

Devices

Synchronization and locking basics

Concurrency

Synchronization primitives

Final workshop

What now?

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate
code

Linked lists

Small allocations

Userspace interface

Syscalls
Filesystems
Devices

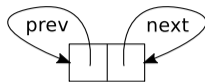
Synchronization and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Linked lists in kernel

- ▶ Doubly-linked
- ▶ Circular



Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code

Linked lists

Small allocations

Userspace interface

Syscalls

Filesystems

Devices

Synchronization
and locking basics

Concurrency

Synchronization
primitives

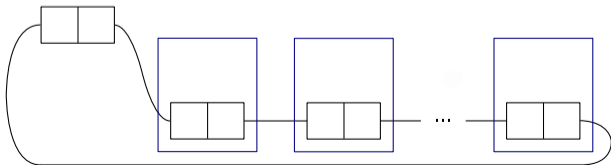
Final workshop

What now?

```
1 struct list_head {  
2     struct list_head *next, *prev;  
3 };
```

Listing 11: Linked list structure

- ▶ Embedded struct list_head
- ▶ *container_of*



Linked lists API (selection)

```
1 /* declare list handle */
2 LIST_HEAD(my_list)
3
4 /* add/delete */
5 list_add(struct list_head *new, struct list_head *head)
6 list_add_tail(struct list_head *new, struct list_head *head)
7 list_del(struct list_head *entry)
8
9 /* query/get elements */
10 int list_empty(const struct list_head *head)
11 list_entry(ptr, type, member)
12 list_first_entry(head, type, member)
13 list_last_entry(head, type, member)
14
15 /* iterate */
16 list_for_each(pos, head)
17 list_for_each_safe(pos, next_tmp, head)
18 list_for_each_entry(pos, head, member)
19 list_for_each_entry_safe(pos, next_tmp, head, member)
20
21 /* see include/linux/list.h for more */
```

Listing 12: Linked lists API (selection)

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code

Linked lists

Small allocations

Userspace interface

Syscalls

Filesystems

Devices

Synchronization
and locking basics

Concurrency

Synchronization
primitives

Final workshop

What now?

Workshop plan

Meeting the Linux kernel

- OS purpose and Linux Hardware, Kernel & toolchain
- (Early) userspace

Writing basic modules

- In-tree vs out-of-tree
- Module boilerplate code
- Linked lists
- Small allocations

Userspace interface

- Syscalls
- Filesystems
- Devices

Synchronization and locking basics

- Concurrency
- Synchronization primitives
- Final workshop

What now?

Small allocations

`kzalloc` allocate (several kB max)

`kfree` free

- ▶ *sizeof* is your friend
- ▶ verify that the return value is `!= NULL`
- ▶ For greater amounts use other methods!²²

```
1 void *kzalloc(size_t size, gfp_t flags) /* flags almost always GFP_KERNEL */  
2 kfree(const void *)
```

Listing 13: Small allocations API

²²e.g. page allocator, slab caches, `vmalloc`

Workshop 4²³

- ▶ Analyze krg.c
- ▶ Fill the blanks as instructed inline
- ▶ Compile and deploy the module
- ▶ Test your module
 - ▶ Try inserting and removing it several times
 - ▶ Try using the 'number' module parameter

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists

Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Workshop 4 solution

Kernel Rookie
Guide Workshop

Andrzej
Pietrasiewicz

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists

Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls

Filesystems

Devices

Synchronization and locking basics

Concurrency

Synchronization primitives

Final workshop

What now?

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls

Filesystems

Devices

Synchronization and locking basics

Concurrency

Synchronization primitives

Final workshop

What now?

Syscalls (1/2)

- ▶ Userspace programs don't 'link' against kernel
- ▶ Syscalls have numbers (e.g. on ARM exit:1, fork:2, read:3, write:4 etc)
- ▶ ABI - defines interaction between binary programs, including convention of passing parameters, receiving results and triggering kernel services²⁴
 - ▶ CPU registers (e.g. on ARM r0, r1, ...)
 - ▶ stack

²⁴Each ABI requires its matching toolchain

Syscalls (2/2)

- ▶ ARM Embedded ABI (here: 32-bit)²⁵
 - ▶ syscall number in r7
 - ▶ syscall arguments in r0...r6
 - ▶ return value in r0
 - ▶ software interrupt #0 triggers syscall execution

Meeting the Linux kernel

OS purpose and Linux Hardware, Kernel & toolchain (Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization primitives
Final workshop

What now?

```
1 # call syscall 'exit' with parameter 0
2 mov r7, #1
3 mov r0, #0
4 svc #0
```

Listing 14: Calling a syscall

- ▶ libc is your friend!
- ▶ strace is your friend, too!

²⁵man syscall, other ARM ABIs exist, e.g. armhf

Workshop 5²⁷

- ▶ Analyze, compile²⁶ and deploy "Hello, world" written in ARM assembly
- ▶ Run it in the virtual machine

²⁶use the provided `compile_hello.sh`

²⁷TODO.txt in W5

Workshop 5 solution

Kernel Rookie
Guide Workshop

Andrzej
Pietrasiewicz

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Adding new syscalls

- ▶ **Don't**
- ▶ Unless you really have to²⁸
- ▶ Or in a training

²⁸Difficult if you want it upstream

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls

Filesystems

Devices

Synchronization and locking basics

Concurrency

Synchronization primitives

Final workshop

What now?

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Architectural overview

- ▶ User
- ▶ Syscall interface
- ▶ VFS (Virtual Filesystem/Virtual Filesystem Switch)
- ▶ 'Real' filesystems (ext2, ext3, ext4....), pseudo filesystems²⁹
- ▶ Block layer
- ▶ Device drivers

²⁹all available ('real' and pseudo) listed in `/proc/filesystems`  56/75

Pseudo filesystems

- ▶ Pseudo filesystems
 - ▶ sysfs
 - ▶ proc
 - ▶ debugfs
 - ▶ configfs
 - ▶ ...
- ▶ No Block layer/Device drivers
- ▶ Particular filesystem driver knows what it means to read/write a 'file'
- ▶ 'Everything is a file' philosophy

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization primitives
Final workshop

What now?

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Devices and drivers

- ▶ One of OS purposes: Manage devices
- ▶ Character devices, block devices
- ▶ entries in `/dev`

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

- ▶ Swiss army knife of all syscalls³⁰
- ▶ Called on device special file file descriptor
- ▶ 'Sub-syscall' number encoded in parameter (interpreted by driver)
- ▶ Optional pointer to user memory

```
1 int ioctl(int fd, unsigned long request, ...);
```

Listing 15: libc's `ioctl()` prototype

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

³⁰man ioctl

Memory access considerations

- ▶ process virtual memory
- ▶ kernel virtual memory
- ▶ accessing user memory from kernel might be non-trivial!³¹
- ▶ `copy_from_user()`, `copy_to_user()`

³¹<https://youtu.be/EWwfMM2AW9g>

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization primitives
Final workshop

What now?

Workshop plan

Meeting the Linux kernel

- OS purpose and Linux Hardware, Kernel & toolchain
- (Early) userspace

Writing basic modules

- In-tree vs out-of-tree
- Module boilerplate code
- Linked lists
- Small allocations

Userspace interface

- Syscalls
- Filesystems
- Devices

Synchronization and locking basics

- Concurrency
- Synchronization primitives
- Final workshop

What now?

Meeting the Linux kernel

- OS purpose and Linux Hardware, Kernel & toolchain
- (Early) userspace

Writing basic modules

- In-tree vs out-of-tree
- Module boilerplate code
- Linked lists
- Small allocations

Userspace interface

- Syscalls
- Filesystems
- Devices

Synchronization and locking basics

- Concurrency**
- Synchronization primitives
- Final workshop

What now?

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization primitives
Final workshop

What now?

Common synchronization primitives 1/3

- ▶ Semaphore³²
 - ▶ Can be incremented and decremented at will
 - ▶ Any thread of execution can do it
 - ▶ If current value 0, next to decrement will be put to sleep
 - ▶ Sleeping until someone increments (first come - first served)
 - ▶ Potentially long waiting time

```
1 /* definition */  
2  
3 /* initial count 1 */  
4 DEFINE_SEMAPHORE(name)  
5  
6 /* any initial count */  
7 struct semaphore my_semaphore = __SEMAPHORE_INITIALIZER(my_semaphore, initial_count)  
8 void sema_init(struct semaphore *sem, int val)  
9  
10 /* increment/decrement */  
11 void up(struct semaphore *sem)  
12 int down_interruptible(struct semaphore *sem)
```

Listing 16: Semaphore API (selection)

³²#include <linux/semaphore.h>

Common synchronization primitives 2/3

▶ Mutex³³

- ▶ Only two states: available and busy ('a binary semaphore')
- ▶ Can only be released by the one who acquired
- ▶ If busy, next to acquire will be put to sleep
- ▶ Sleeping until someone releases (first come - first served)
- ▶ Avoid long waiting time

```
1 /* definition */  
2 DEFINE_MUTEX(mutexname)  
3 mutex_init(mutex)  
4  
5 /* acquire/release */  
6 void mutex_lock(struct mutex *lock)  
7 void mutex_unlock(struct mutex *lock)  
8  
9 /* acquire or return, without waiting */  
10 int mutex_trylock(struct mutex *lock)
```

Listing 17: Mutex API (selection)

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

³³`#include <linux/mutex.h>`

Common synchronization primitives 3/3

▶ Spinlock³⁴

- ▶ Busy waiting (can be good!)
- ▶ While contended, no guarantee of 'fair' acquiring across all architectures
- ▶ Sometimes turning off interrupts on local CPU required
- ▶ Must be used if sleeping prohibited in a context
- ▶ Supposed to be kept for very short time

```
1  /* initialize */  
2  DEFINE_SPINLOCK(lock)  
3  
4  /* acquire/release */  
5  void spin_lock(spinlock_t *lock)  
6  void spin_unlock(spinlock_t *lock)  
7  
8  spin_lock_irqsave(lock, flags)  
9  void spin_unlock_irqrestore(spinlock_t *lock, unsigned long flags)  
10  
11 /* acquire or return, without waiting */  
12 int spin_trylock(spinlock_t *lock)
```

Listing 18: Spinlock API (selection)

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

³⁴`#include <linux/spinlock.h>`

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

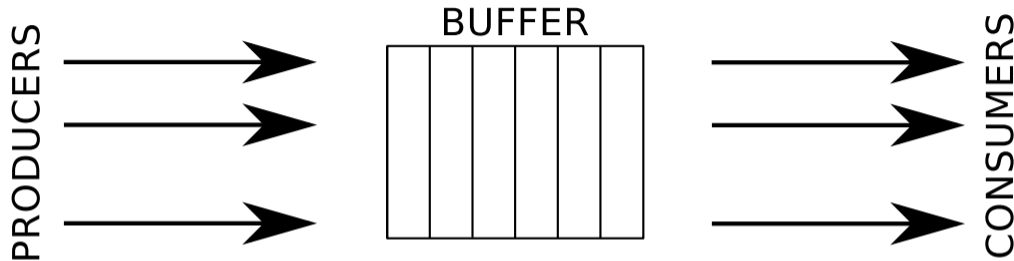
Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization primitives
Final workshop

What now?

Producer and consumer (aka Bounded buffer) problem



Meeting the Linux kernel

OS purpose and Linux Hardware, Kernel & toolchain (Early) userspace

Writing basic modules

In-tree vs out-of-tree Module boilerplate code Linked lists Small allocations

Userspace interface

Syscalls Filesystems Devices

Synchronization and locking basics

Concurrency Synchronization primitives Final workshop

What now?

```
1 down(empty_count) /* at init: buffer size */
2 lock(buffer)
3 put_into(buffer)
4 unlock(buffer)
5 up(fill_count)
```

Listing 19: Producer pseudo code

```
1 down(fill_count) /* at init: 0 */
2 lock(buffer)
3 get_from(buffer)
4 unlock(buffer)
5 up(empty_count)
```

Listing 20: Consumer pseudo code

- ▶ Analyze the userspace program (syscall.c) and discover what it does
- ▶ Analyze arch/arm/kernel/sys_arm.c in kernel sources and fill the blanks as instructed inline
- ▶ Rebuild the kernel
- ▶ Try running ./syscall in the virtual machine
 - ▶ ./syscall -p <number> to put element into buffer
 - ▶ ./syscall -g to get element from buffer
 - ▶ You might need to put either into background with & (e.g. while getting from empty buffer)

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Workshop 6 solution

Kernel Rookie
Guide Workshop

Andrzej
Pietrasiewicz

Meeting the Linux
kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic
modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization
and locking basics

Concurrency
Synchronization
primitives

Final workshop

What now?

Workshop plan

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel & toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization primitives
Final workshop

What now?

Meeting the Linux kernel

OS purpose and Linux
Hardware, Kernel &
toolchain
(Early) userspace

Writing basic modules

In-tree vs out-of-tree
Module boilerplate
code
Linked lists
Small allocations

Userspace interface

Syscalls
Filesystems
Devices

Synchronization and locking basics

Concurrency
Synchronization
primitives
Final workshop

What now?

Appendix

```
1 #!/bin/bash
2
3 LINUX_SOURCE=../linux-4.18
4 INITRAMFS_PATH=${PWD}/initramfs
5 MODULES_PATH=${PWD}/..
6 INITRAMFS_LIST_PATH=${PWD}/initramfs
7
8 cd ${LINUX_SOURCE}
9 scripts/gen_initramfs_list.sh -o ${INITRAMFS_PATH}/initramfs.cpio \
10  ${INITRAMFS_LIST_PATH}/initramfs.list
```

Listing 21: generate_initramfs.sh

```
1 #!/bin/bash
2
3 LINUX_SOURCE=../linux-4.18
4 INITRAMFS_PATH=initramfs
5 QEMU_PREFIX=
6
7 ${QEMU_PREFIX}qemu-system-arm \
8   -nographic \
9   -M vexpress-a9 \
10  -kernel ${LINUX_SOURCE}/arch/arm/boot/zImage \
11  -append "root=/dev/ram rdinit=/busybox-armv7r init console=ttyAMA0,115200n8 earlyprintk" \
12  -dtb ${LINUX_SOURCE}/arch/arm/boot/dts/vexpress-v2p-ca9.dtb \
13  -initrd ${INITRAMFS_PATH}/initramfs.cpio
```

Listing 22: vexpress_run.sh

```
1 file /busybox-armv7r ../workshops/initramfs/busybox-armv7r 0755 0 0
2 dir /sbin 0755 0 0
3 dir /proc 0755 0 0
4 dir /sys 0755 0 0
5 dir /bin 0755 0 0
6 dir /usr 0755 0 0
7 dir /usr/sbin 0755 0 0
8 dir /usr/bin 0755 0 0
9 dir /etc 0755 0 0
10 dir /etc/init.d 0755 0 0
11 file /etc/inittab ../workshops/initramfs/inittab 0755 0 0
12 file /etc/init.d/rcS ../workshops/initramfs/rcS 0755 0 0
13 file /syscall ../workshops/initramfs/syscall 0755 0 0
14 file /strace ../workshops/initramfs/strace 0755 0 0
```

Listing 23: initramfs.list

```
1 ::sysinit:/etc/init.d/rcS
2 ::askfirst:-/bin/sh
3 ::ctrlaltdel:/sbin/reboot
4 ::shutdown:/sbin/swapoff -a
5 ::shutdown:/sbin/umount -a -r
6 ::restart:/sbin/init
```

Listing 24: inittab

```
1 #!/busybox-armv7r sh
2
3 /busybox-armv7r --install
4
5 mount -t proc none /proc
6 mount -t sysfs none /sys
7
8 mdev -s
9 echo /sbin/mdev > /proc/sys/kernel/hotplug
```

Listing 25: rcS

```
1 From 3aae8f646785afab58f2195111d6280029af7ac0 Mon Sep 17 00:00:00 2001
2 From: Andrzej Pietrasiewicz <andrzej.p@samsung.com>
3 Date: Fri, 17 Aug 2018 13:58:36 +0200
4 Subject: [PATCH] kbuild: make sorting initramfs contents independent of locale
5
6 Some LANG values (e.g. pl_PL.UTF-8) cause the sort command to output
7 files before their parent directories, which makes them inaccessible for
8 the kernel. In other words, when the kernel populates the rootfs, it is
9 unable to create files whose parent directories have not been yet created.
10
11 This patch makes sorting use the default (LANG=C) locale, which results in
12 correctly laid out initramfs images (parent directories before files).
13
14 Signed-off-by: Andrzej Pietrasiewicz <andrzej.p@samsung.com>
15 ---
16 scripts/gen_initramfs_list.sh | 2 +-
17 1 file changed, 1 insertion(+), 1 deletion(-)
18
19 diff --git a/scripts/gen_initramfs_list.sh b/scripts/gen_initramfs_list.sh
20 index 10e528b..0aad760 100755
21 --- a/scripts/gen_initramfs_list.sh
22 +++ b/scripts/gen_initramfs_list.sh
23 @@ -174,7 +174,7 @@ dir_filelist() {
24     ${dep_list}header "$1"
25
26     srcdir=$(echo "$1" | sed -e 's://*/:/:g')
27 -   dirlist=$(find "${srcdir}" -printf "%p_%m_%U_%G\n" | sort)
28 +   dirlist=$(find "${srcdir}" -printf "%p_%m_%U_%G\n" | LANG=C sort)
29
30     # If $dirlist is only one line, then the directory is empty
31     if [ "$(echo "${dirlist}" | wc -l)" -gt 1 ]; then
32     --
33 2.7.4
```



```
1 From 21cf89b5a128b6a3390efd6c160f67edfde20fa9 Mon Sep 17 00:00:00 2001
2 From: Bartlomiej Zolnierkiewicz <b.zolnierkie@samsung.com>
3 Date: Fri, 31 Aug 2018 11:53:36 +0200
4 Subject: [PATCH] Revert "irqdomain: Don't set type when mapping an IRQ"
5
6 This reverts commit 1e2a7d78499ec8859d2b469051b7b80bad3b08aa.
7
8 This makes ARM QEMU vexpress platform work again.
9
10 Signed-off-by: Bartlomiej Zolnierkiewicz <b.zolnierkie@samsung.com>
11 ---
12 include/linux/irqdomain.h | 3 ---
13 kernel/irq/irqdomain.c    | 23 +++++-----
14 2 files changed, 5 insertions(+), 21 deletions(-)
15
16 diff --git a/include/linux/irqdomain.h b/include/linux/irqdomain.h
17 index dccfa65..e93742d 100644
18 --- a/include/linux/irqdomain.h
19 +++ b/include/linux/irqdomain.h
20 @@ -526,9 +526,6 @@ static inline int irq_domain_alloc_irqs(struct irq_domain *domain,
21     return -1;
22 }
23
24 -static inline void irq_domain_free_irqs(unsigned int virq,
25 -    unsigned int nr_irqs) { }
26 -
27 static inline bool irq_domain_is_hierarchy(struct irq_domain *domain)
28 {
29     return false;
```

Listing 27: irqdomain patch 1/3

```
1 diff --git a/kernel/irq/irqdomain.c b/kernel/irq/irqdomain.c
2 index 5d9fc01b..75f175b 100644
3 --- a/kernel/irq/irqdomain.c
4 +++ b/kernel/irq/irqdomain.c
5 @@ -744,7 +744,6 @@ static void of_phandle_args_to_fwspec(struct of_phandle_args *irq_data,
6  unsigned int irq_create_fwspec_mapping(struct irq_fwspec *fwspec)
7  {
8  struct irq_domain *domain;
9  - struct irq_data *irq_data;
10  irq_hw_number_t hwirq;
11  unsigned int type = IRQ_TYPE_NONE;
12  int virq;
13 @@ -792,11 +791,7 @@ unsigned int irq_create_fwspec_mapping(struct irq_fwspec *fwspec)
14  * it now and return the interrupt number.
15  */
16  if (irq_get_trigger_type(virq) == IRQ_TYPE_NONE) {
17  - irq_data = irq_get_irq_data(virq);
18  - if (!irq_data)
19  - return 0;
20  -
21  irqd_set_trigger_type(irq_data, type);
22  + irq_set_irq_type(virq, type);
23  return virq;
24  }
25
26 @@ -816,18 +811,10 @@ unsigned int irq_create_fwspec_mapping(struct irq_fwspec *fwspec)
27  return virq;
28  }
29
30 - irq_data = irq_get_irq_data(virq);
31 - if (!irq_data) {
32 - if (irq_domain_is_hierarchy(domain))
33 - irq_domain_free_irqs(virq, 1);
```

```
1 -     else
2 -         irq_dispose_mapping(virq);
3 -     return 0;
4 - }
5 -
6 - /* Store trigger type */
7 - irqd_set_trigger_type(irq_data, type);
8 -
9 - /* Set type if specified and different than the current one */
10 + if (type != IRQ_TYPE_NONE &&
11 +     type != irq_get_trigger_type(virq))
12 +     irq_set_irq_type(virq, type);
13     return virq;
14 }
15 EXPORT_SYMBOL_GPL(irq_create_fwspec_mapping);
16 --
17 1.9.1
```

Listing 29: irqdomain patch 3/3

```
1 From cb1d564a98486d331f9fb164e8d266d9aedd0cd8 Mon Sep 17 00:00:00 2001
2 From: Andrzej Pietrasiewicz <andrzej.p@samsung.com>
3 Date: Fri, 5 Oct 2018 10:08:49 +0200
4 Subject: [PATCH] krg locking workshop
5
6 Signed-off-by: Andrzej Pietrasiewicz <andrzej.p@samsung.com>
7 ---
8 arch/arm/kernel/sys_arm.c | 116 +++++
9 arch/arm/tools/syscall.tbl |  2 +
10 2 files changed, 118 insertions(+)
11
12 diff --git a/arch/arm/kernel/sys_arm.c b/arch/arm/kernel/sys_arm.c
13 index bdf7514..072cc11 100644
14 --- a/arch/arm/kernel/sys_arm.c
15 +++ b/arch/arm/kernel/sys_arm.c
16 @@ -27,6 +27,9 @@
17  #include <linux/ipc.h>
18  #include <linux/uaccess.h>
19  #include <linux/slab.h>
20  +#include <linux/semaphore.h>
21  +#include <linux/mutex.h>
22  +#include <linux/list.h>
23
24  /*
25   * Since loff_t is a 64 bit type we avoid a lot of ABI hassle
26  @@ -37,3 +40,116 @@ asmlinkage long sys_arm_fadvise64_64(int fd, int advice,
27  {
28     return ksys_fadvise64_64(fd, offset, len, advice);
29 }
30 +
31 +/***** Kernel Rookie Guide OSSEU'18 Edinburgh *****/
32  #define KRG_BUF_SIZE      13
33  +
```

```
1 +struct krg_job {
2 + struct list_head entry;
3 + int id;
4 +};
5 +
6 +static struct semaphore krg_empty_count = __SEMAPHORE_INITIALIZER(krg_empty_count, KRG_BUF_SIZE);
7 +static struct semaphore krg_fill_count = __SEMAPHORE_INITIALIZER(krg_fill_count, 0);
8 +
9 +static LIST_HEAD(krg_jobs);
10 +static DEFINE_MUTEX(krg_mutex);
11 +
12 +SYSCALL_DEFINE1(krg_put, int, id)
13 +{
14 + struct krg_job *job;
15 +
16 + job = NULL;
17 + /*
18 +  * TODO: allocate struct krg_job
19 +  * use kzalloc()
20 +  */
21 + if (!job)
22 + return PTR_ERR(job);
23 +
24 + /* TODO: set job id */
25 +
26 + /*
27 +  * TODO: decrement empty count, use
28 +  *
29 +  * if (down_interruptible()) {
30 +  * TODO: avoid memory leak
31 +  * use kfree()
32 +  *
33 +  * // upper layers can handle retrying
```

Listing 31: krg workshop patch 2/5

```
1 + * return -ERESTARTSYS;
2 + * }
3 + */
4 +
5 + /*
6 + * TODO: ensure only one process accesses the buffer
7 + * use mutex_lock()
8 + */
9 +
10 + printk(">_put:%d\n", job->id);
11 +
12 + /*
13 + * TODO: add job to the list
14 + * use list_add_tail()
15 + */
16 +
17 + /*
18 + * TODO: let other processes access the buffer
19 + * use mutex_unlock()
20 + */
21 +
22 + /*
23 + * TODO: increment fill count
24 + * use up()
25 + */
26 +
27 + return 0;
28 +}
29 +
30 +SYSCALL_DEFINE0(krg_get)
31 +{
32 + struct krg_job *job;
33 + int ret;
```

Listing 32: krg workshop patch 3/5

```
1 +
2 + /*
3 +  * TODO: decrement fill count, use
4 +  * if (down_interruptible())
5 +  * return -ERESTARTSYS;
6 +  */
7 +
8 + /*
9 +  * TODO: ensure only one process accesses the buffer
10 +  * use mutex_lock()
11 +  */
12 +
13 + /*
14 +  * TODO: let 'job' point to the first element of the buffer
15 +  * use job = list_first_entry()
16 +  */
17 +
18 + printk("<_get:%d\n", job->id);
19 +
20 + ret = job->id;
21 +
22 + /*
23 +  * TODO: remove the element from the buffer
24 +  * use list_del()
25 +  */
26 +
27 + /*
28 +  * TODO: avoid memory leak
29 +  * use kfree()
30 +  */
31 +
32 + /*
33 +  * TODO: let other processes access the buffer
```

Listing 33: krg workshop patch 4/5

```
1 + * use mutex_unlock()
2 + */
3 +
4 + /*
5 + * TODO: increment empty count
6 + * use up()
7 + */
8 +
9 + return ret;
10 +}
11 diff --git a/arch/arm/tools/syscall.tbl b/arch/arm/tools/syscall.tbl
12 index fbc74b5..febf960 100644
13 --- a/arch/arm/tools/syscall.tbl
14 +++ b/arch/arm/tools/syscall.tbl
15 @@ -413,3 +413,5 @@
16 396 common pkey_free sys_pkey_free
17 397 common statx sys_statx
18 398 common rseq sys_rseq
19 +399 eabi krg_put sys_krg_put
20 +400 eabi krg_get sys_krg_get
21 --
22 2.7.4
```

Listing 34: krg workshop patch 5/5


```
1 1. Find container of definition in 4.18 and try understanding it
2
3 2. If more than one found choose the one whose location looks most general
4
5 3. https://elixir.bootlin.com, use 4.18
```

Listing 35: W0/TODO.txt

```
1 1. Patch the kernel with the provided patches
2
3 In kernel source directory:
4
5 patch -p1 < PATH_TO_WORKSHOPS/kernel-patches/0001-kbuild....
6 patch -p1 < PATH_TO_WORKSHOPS/kernel-patches/0001-Revert....
7 patch -p1 < PATH_TO_WORKSHOPS/kernel-patches/0001-krng....
8
9 2. Configure & cross-compile the kernel for vexpress
10
11 In kernel source directory:
12
13 make ARCH=arm vexpress_defconfig
14 make -j8 ARCH=arm CROSS_COMPILE=arm-linux-gnueabi-
15
16 3. Analyze the vexpress-run.sh script
17
18 4. Run emulated machine in qemu with the script
19 ^A-X to exit qemu
```

Listing 36: W1/TODO.txt

```
1 1. Add a short text file to initramfs image
2
3 Use initramfs.list - pay special attention to host paths which are relative to
4 kernel source toplevel directory!
5
6 2. Use generate_initramfs.sh to build initramfs
7
8 3. Run the virtual machine and verify new file's presence and contents
```

Listing 37: W2/TODO.txt

```
1 1. Write, compile and deploy (initramfs!) a module which does (almost) nothing
2
3 Module's init shall output a welcome string
4 Module's exit shall output a goodbye string
5
6 Use the W3 directory and a Makefile there.
7
8 Remember about ARCH=arm and CROSS_COMPILE=arm-linux-gnueabi-
9
10 2. Load/unload the module in the virtual machine
```

Listing 38: W3/TODO.txt

```
1 ifneq (${KERNELRELEASE},)
2
3 obj-m := hello.o
4
5 else
6
7 # modify appropriately
8 LINUX_SOURCE := ../../linux-4.18
9
10 PWD := $(shell pwd)
11
12 default:
13     ${MAKE} -C ${LINUX_SOURCE} SUBDIRS=${PWD} modules
14
15 clean:
16     ${MAKE} -C ${LINUX_SOURCE} SUBDIRS=${PWD} clean
17
18 endif
```

Listing 39: W3/Makefile

```
1 1. Analyze krg.c
2
3 2. Fill the blanks as instructed inline
4
5 3. Compile and deploy the module
6
7 Remember about ARCH=arm and CROSS_COMPILE=arm-linux-gnueabi-
8
9 4. Test your module
10     Try inserting and removing it several times
11     Try using the 'number' module parameter
```

Listing 40: W4/TODO.txt

```
1 ifneq ( ${KERNELRELEASE}, )
2
3 obj-m := krg.o
4
5 else
6
7 # modify appropriately
8 LINUX_SOURCE := ../../linux-4.18
9
10 PWD := $(shell pwd)
11
12 default:
13     ${MAKE} -C ${LINUX_SOURCE} SUBDIRS=${PWD} modules
14
15 clean:
16     ${MAKE} -C ${LINUX_SOURCE} SUBDIRS=${PWD} clean
17
18 endif
```

Listing 41: W4/Makefile

```
1 /***** Kernel Rookie Guide OSSEU'18 Edinburgh *****/
2 *
3 * (C) Andrzej Pietrasiewicz <andrzej.p@samsung.com>
4 *
5 * This program is free software; you can redistribute it and/or modify
6 * it under the terms of the GNU General Public License version 2 as
7 * published by the Free Software Foundation.
8 *
9 */
10 #include <linux/module.h>
11 #include <linux/list.h>
12 #include <linux/slab.h>
13
14 struct data {
15     int id;
16     struct list_head entry;
17 };
18
19 static LIST_HEAD(collection);
20
21 /*
22  * The module allocates 'number' instances (default: 3) of struct data
23  * and adds them to 'collection'.
24  */
25 static int number = 3;
26 module_param(number, int, 0);
27
28 int krg_init(void)
29 {
30     struct data *data;
31     int i;
32
33     /* we want our list to have some length, no more than 1000 */
```

```
1  if (number <= 0 || number > 1000)
2      return -EINVAL;
3
4  /* TODO: ensure appropriate 'number' of iterations */
5  for (i = 0; i < 0; ++i) {
6      /*
7       * TODO: allocate a struct data
8       * use kzalloc()
9       */
10     data = NULL;
11     if (!data)
12         goto no_memory;
13
14     /* TODO: set the id member of struct data to i */
15
16     /*
17      * TODO: add struct data after the last element of collection
18      * use list_add_tail()
19      */
20
21     printk("%s added: %d %px\n", __func__, data->id, data);
22 }
23
24 return 0;
25
26 no_memory:
27 while (i--) {
28     /*
29      * TODO: let 'data' point to last element of collection
30      * use list_last_entry()
31      */
32
33     /*
```

```
1      * TODO: remove the element from the collection
2      * use list_del()
3      */
4
5     /*
6     * TODO: free the memory occupied by the removed element
7     * use kfree()
8     */
9 }
10 return -ENOMEM;
11 }
12 module_init(krg_init);
13
14 void krg_exit(void)
15 {
16     struct data *data;
17
18     while (number-- > 0) {
19         /*
20         * TODO: let 'data' point to last element of collection
21         * use list_last_entry()
22         */
23
24         /*
25         * TODO: remove the element from the collection
26         * use list_del()
27         */
28
29         printk("%s removed: %d@%px\n", __func__, data->id, data);
30
31         /*
32         * TODO: free the memory occupied by the removed element
33         * use kfree()

```

```
1     */
2   }
3 }
4 module_exit(krg_exit);
5
6 MODULE_LICENSE("GPL");
```

Listing 45: W4/krg.c 4/4

```
1 1. Analyze, compile and deploy 'Hello, world' written in ARM assembly
2
3 See hello.S
4
5 use the provided compile_hello.sh
6
7 __NR_write
8 __NR_exit
9
10 symbolic numbers of write and exit syscalls
11
12 necessary ARM assembly:
13
14 mov rX, <immediate> - store <immediate> value in register rX
15 ldr rX, <address> - store <address> value in register rX
16 svc <immediate> - trigger software interrupt number <immediate>
17
18 2. Run it in the virtual machine
```

Listing 46: W5/TODO.txt


```
1 #!/bin/bash
2
3 LINUX_SOURCE=../../linux-4.18
4 #
5 # IMPORTANT!
6 #
7 # run these commands once before running this script:
8 #
9 # cd ${LINUX_SOURCE}
10 # make ARCH=arm headers_install
11 #
12 arm-linux-gnueabi-gcc hello.S -c -static -I${LINUX_SOURCE}/usr/include
13 arm-linux-gnueabi-ld hello.o -o hello
```

Listing 47: W5/compile_hello.sh

```
1 #include "linux/unistd.h"
2 .data
3 hello_string:
4     .ascii "Hello_\World\n"
5 .text
6 .globl _start
7 _start:
8     mov r7, #__NR_write
9     mov r0, #1
10    ldr r1,=hello_string
11    mov r2, #12
12    svc #0
13
14    mov r7, #__NR_exit
15    mov r0, #0
16    svc #0
```

```
1 1. Analyze the userspace program (syscall.c) and discover what it does
2
3 2. Analyze arch/arm/kernel/sys_arm.c in kernel sources and fill the blanks
4 as instructed inline
5
6 3. Rebuild the kernel
7
8 Remember about ARCH=arm CROSS_COMPILE=arm-linux-gnueabi-
9
10 4. Try running ./syscall in the virtual machine
11
12     ./syscall -p <number> to put element into buffer
13     ./syscall -g to get element from buffer
14     You might need to put either into background with &
15     (e.g. while getting from empty buffer)
16
17 When trying to get from empty buffer the process will sleep waiting
18 on the semaphore. If you don't put the process into background you
19 have no other shell available to put the element into buffer.
20 Converse situation happens when trying to put into full buffer.
```

Listing 49: W6/TODO.txt

```
1 #include <ctype.h>
2 #include <errno.h>
3 #include <getopt.h>
4 #include <stdbool.h>
5 #include <stdio.h>
6 #include <stdlib.h>
7 #include <unistd.h>
8 #include "linux/unistd.h"
9
10 int main(int argc, char *argv[])
11 {
12     int value;
13
14     while ((value = getopt(argc, argv, "p:g")) != -1)
15         switch (value) {
16             case 'p':
17                 if (!isdigit(optarg[0])) {
18                     fprintf(stderr, "Cannot parse put argument!\n");
19                     exit(EXIT_FAILURE);
20                 }
21                 value = atoi(optarg);
22                 return (syscall(__NR_krg_put, value) < 0) ? errno : 0;
23
24             case 'g':
25                 value = syscall(__NR_krg_get);
26                 return (value < 0) ? errno : 0;
27
28             default:
29                 exit(EXIT_FAILURE);
30         }
31
32     exit(EXIT_SUCCESS);
33 }
```