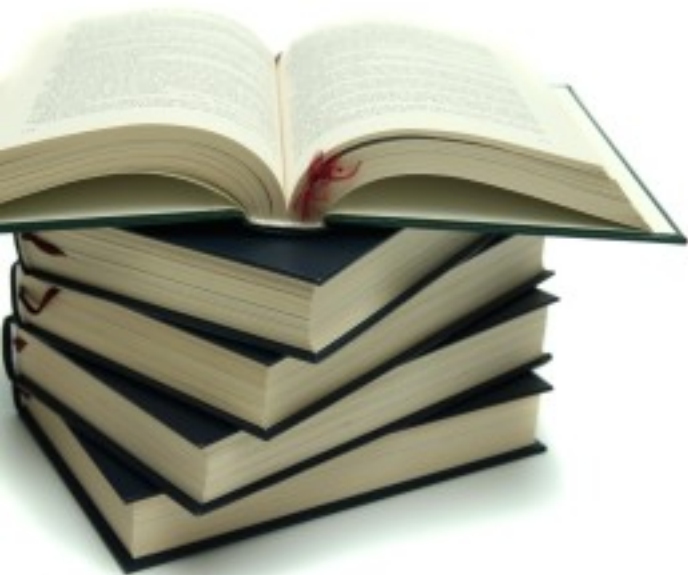


- Ce operație se face în cadrul rutinei de tratare pentru o întrerupere de pachet trimis?
- Care este numărul minim de pași necesar pentru a ruta o adresă IPv4 prin ruta default, dacă adresa nu este în cache-ul de rute și avem cel puțin o rută pentru fiecare din măștile de rețea posibile ?
- Ce operație skb se folosește în procesul de încapsulare al unui pachet? Dar pentru decapsulare?



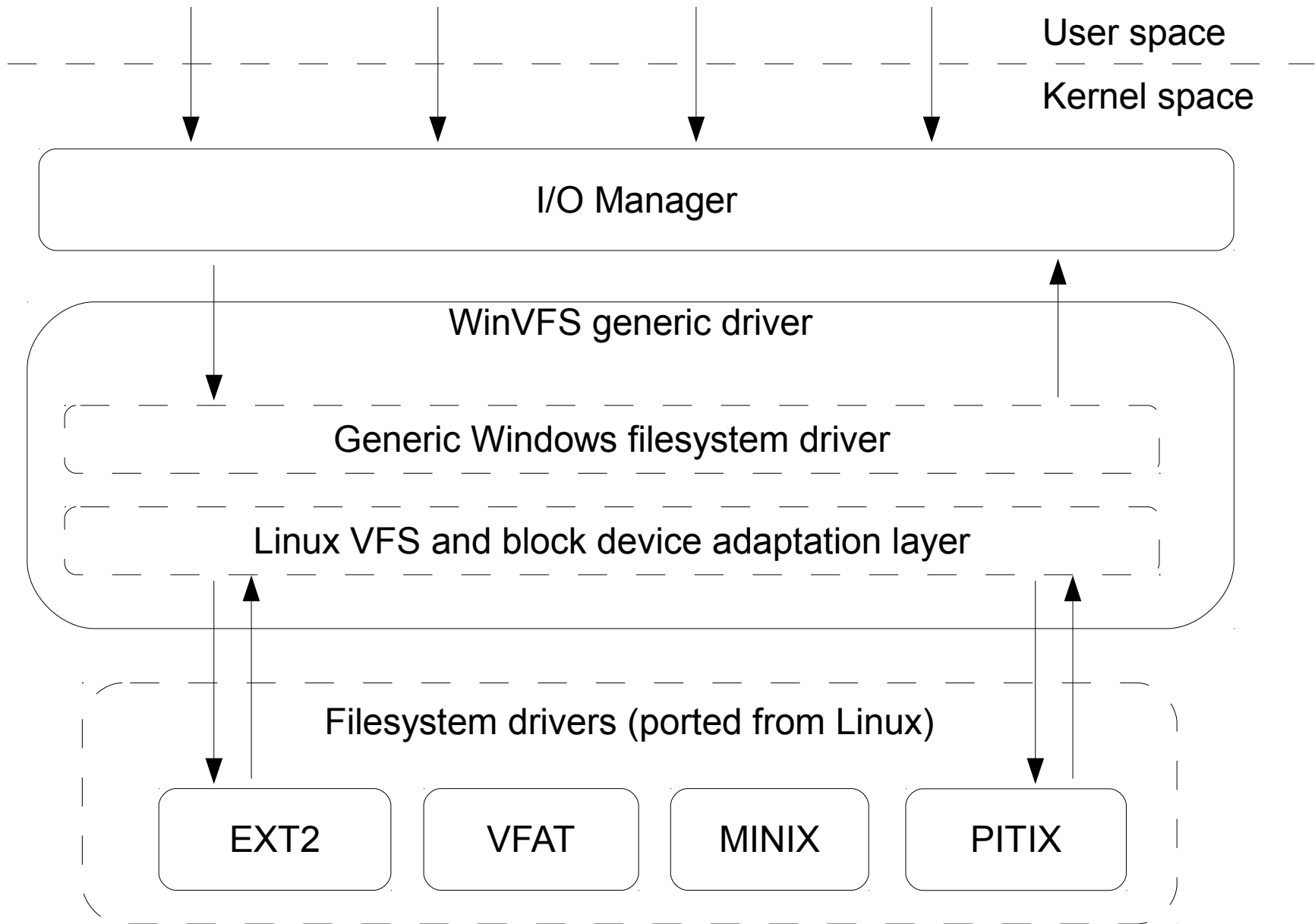
10

Linux kernel library

6 mai 2010

- Initially started
 - As a revival of the WinVFS project but in such a way that we could keep up with the Linux change rate
 - To test a new idea: FTP server as a portable way of offering access to Linux filesystems
- Ended up
 - An infrastructure which allows one to reuse generic Linux kernel code
- Related areas
 - UML
 - CoLinux
 - FuSE, Ndiswrapper
 - Paravirtualization

- Create an Windows ext2 driver as well drivers for other Linux filesystems (reiserfs)
- Completely reuse Linux filesystem drivers: just recompile the driver source code
- WinVFS = the infrastructure needed for complete code reutilization



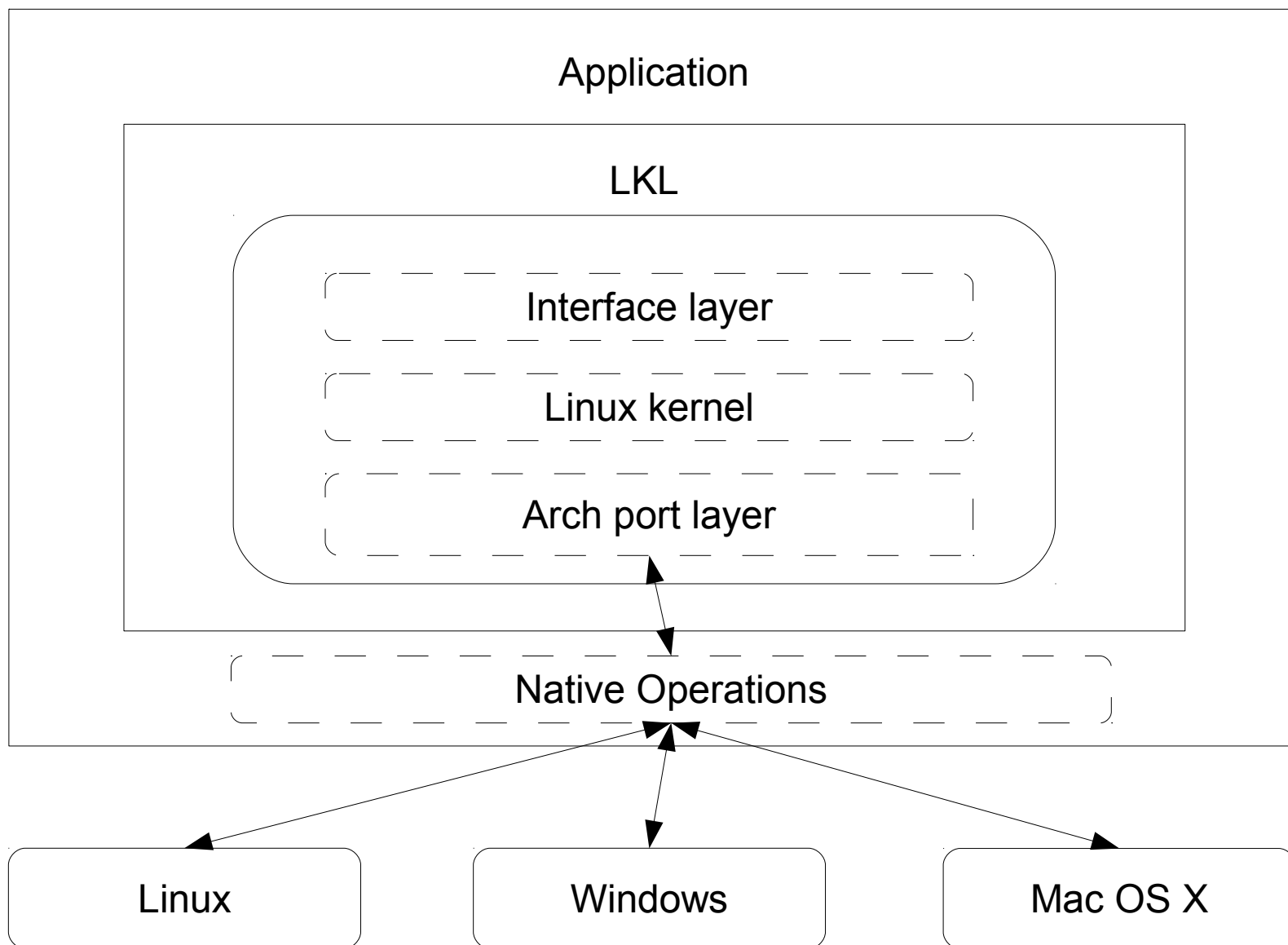
- Generic filesystem driver:
 - Read-only support only
- Adaptation layers:
 - partial porting, partial reimplementaion of the VFS primitives needed by drivers
 - A lot of the generic Linux code was pulled in because of VFS dependencies on various subsystems
- Drivers ported: ext2, minix, vfat
 - Trivial source code modification required (compiler related)

bitops.h config.h errno.h fd.h init.h minix_fs.h msdos_fs.h quota.h
slab.h time.h blkdev.h ctype.h ext2_fs.h file.h ioctl.h minix_fs_i.h
msdos_fs_i.h quotaops.h smp_lock.h types.h blk.h ext2_fs_i.h
fs.h kdev_t.h minix_fs_sb.h msdos_fs_sb.h rwsem.h spinlock.h
wait.h byteorder dcache.h ext2_fs_sb.h fs_struct.h kernel.h
mm.h nls.h rwsem-spinlock.h stat.h capability.h dirent.h
fat_cvf.h highmem.h list.h module.h pagemap.h sched.h
stddef.h compiler.h dnotify.h fcntl.h highuid.h locks.h mount.h
posix_types.h semaphore.h string.h ./mm/page_alloc.c
./mm/kmem_cache.c ./mm/filemap.c ./fs/inode.c ./fs/file_table.c
./fs/attr.c ./fs/namespace.c ./fs/bad_inode.c ./fs/dcache.c
./fs/namei.c ./fs/buffer.c ./fs/readdir.c ./fs/open.c ./fs/super.c ./fs/block_dev.c
./fs/read_write.c ./fs/devices.c ./lib/vsprintf.c ./lib/string.c ./lib/ctype.c

- Switched to mingw
- Switched to 2.6 kernel
- TotalCommander plugins
- Security attributes: Linux – Windows adaptation
- We proved it is possible to completely reuse Linux filesystem drivers code to create Windows drivers
- Switching to 2.6 posed significant challenges
- Keeping track with 2.6 development became impractical

- Allow applications to reuse code from the Linux kernel without needing to separate, isolate and extract it from Linux
- Run in as diverse environments as possible: cross OS, cross platforms, both kernel and user
- Allow full Linux subsystem to be reuse (e.g. filesystem drivers, TCP/IP stack)
- Linux kernel modifications should be isolated (for easy tracking of Linux kernel development)
- Easy to use (from application point of view)

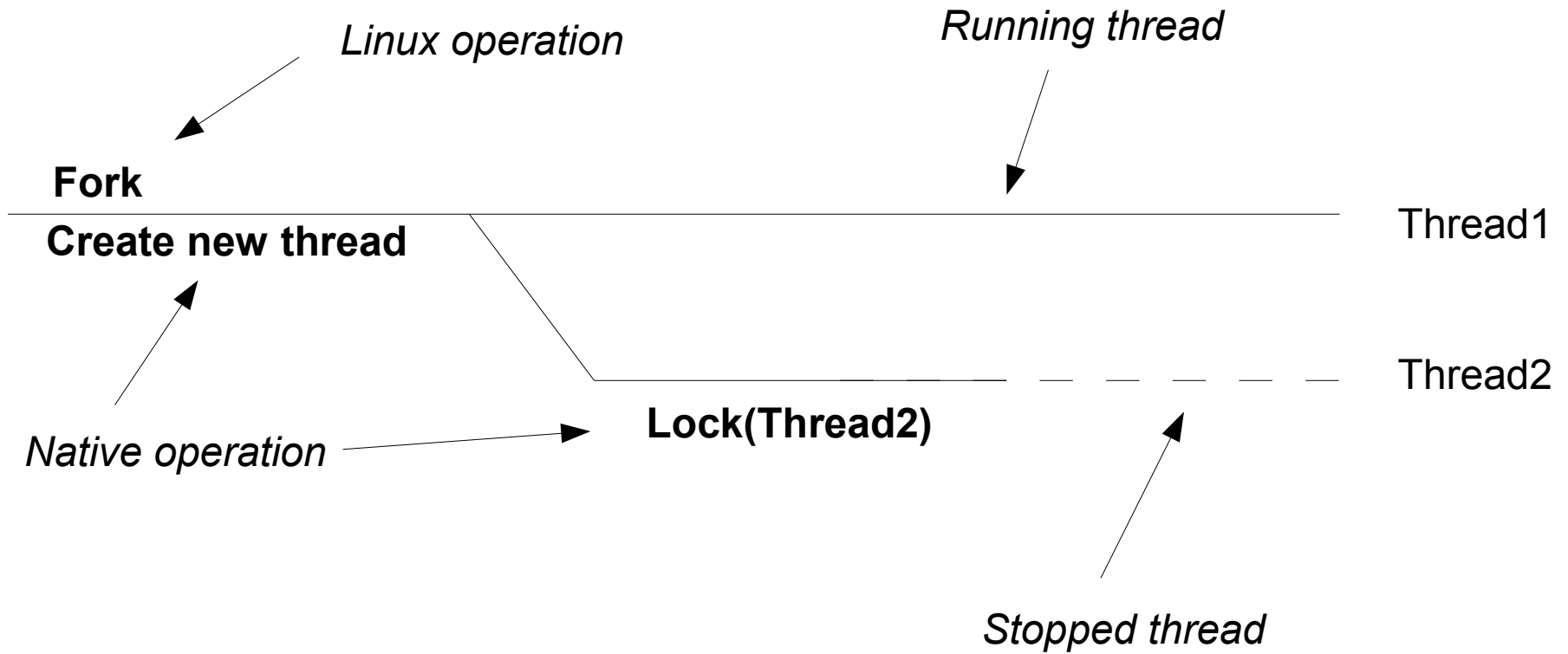
- Make it a library
- The library should contain the full Linux kernel
- Highly customizable – make menuconfig
- Implement it as a new arch port layer
- API based on the Linux system call interface
- Offload some operations to application
- No user / kernel separation or abstractions

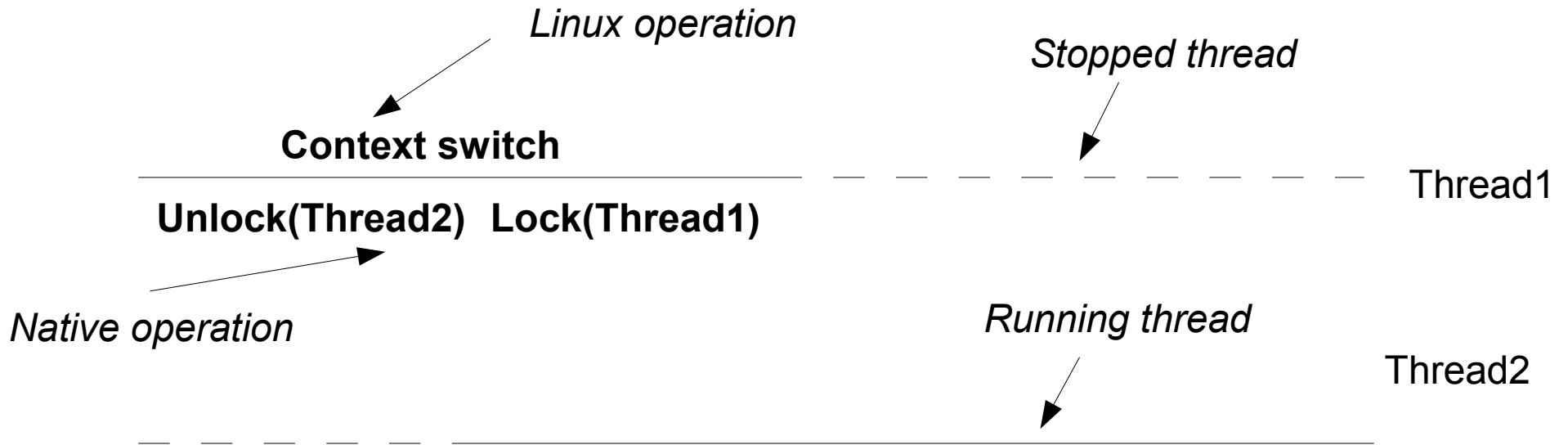


- Offers services needed by the Linux kernel (e.g. memory management, thread management, time management, etc.)
- By design, the operations are basic and as generic as possible
- It is the role of the arch port layer to map these operations to the services required by the Linux kernel

- Ikl arch is a “non-MMU” arch
- “Physical memory” allocated by the native environment
- Initially: allocate the whole physical memory during initialization
- Later: use native operations to allocate memory
 - Hot plug memory

- No need for user processes, but...
- We need to support kernel threads
- Micro/internal LKL threading – discarded
- Support from the execution environment
- Put the Linux scheduler in control:
 - Each thread has a control semaphore
 - Native operations for semaphore control





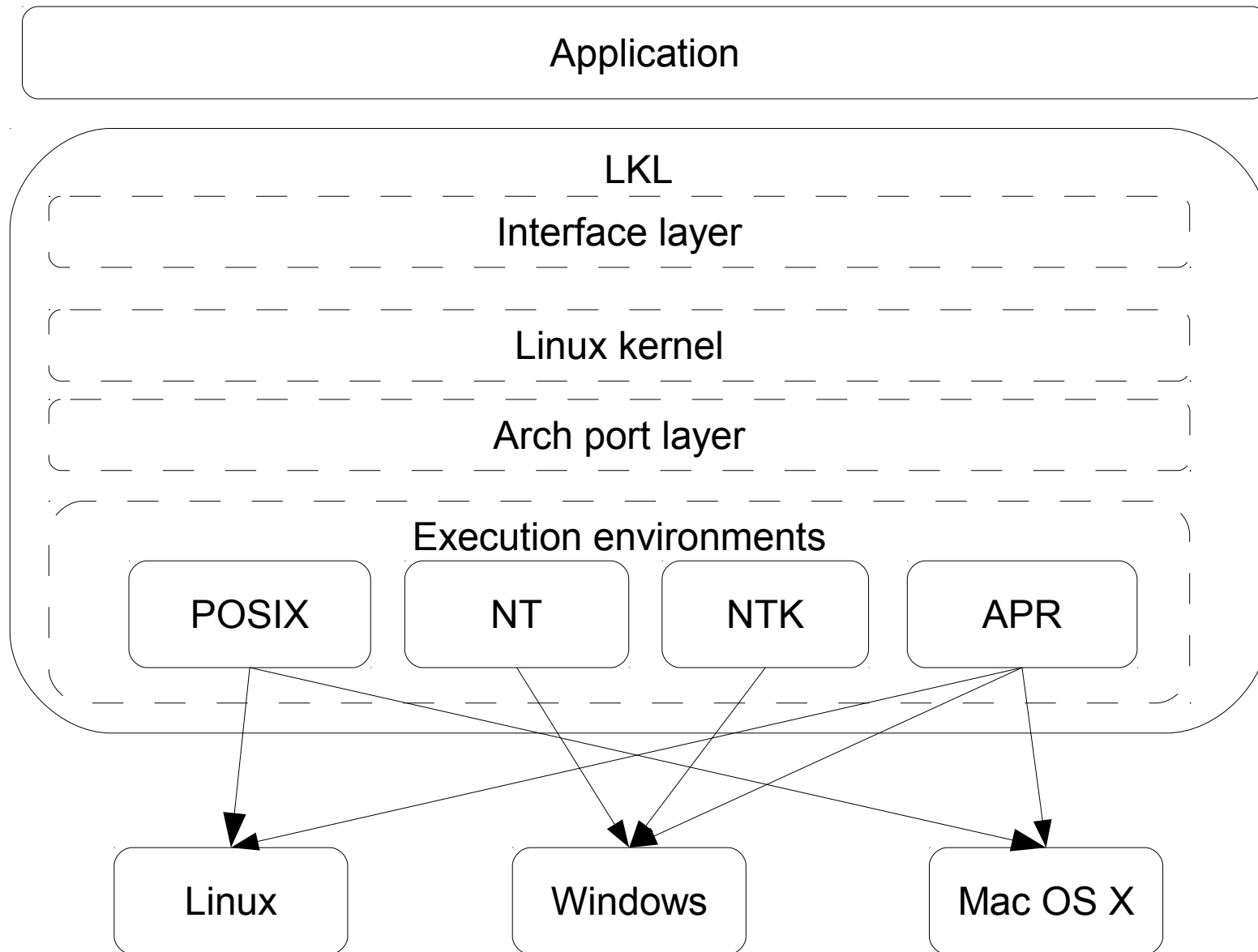
- LKL needs drivers to interact with the exterior
 - Native part - “the hardware”
 - Linux part – a Linux device driver
- How do we communicate between the two parts?
 - Linux -> Native: direct function calls
 - Native -> Linux: “interrupts”
- Why interrupts?
 - The simplest way of running Linux code in the proper context

- Disc driver
 - “Hardware” = file, partition
 - “Hardware” = device object
- Network driver
 - “Hardware” = interface
 - “Hardware” = socket
- Timer driver
- Console driver

- The application can trigger IRQs
- The Linux kernel will pick it up and run the associated interrupt handler
- LKL does not support SMP
- We need to serialize the interrupt handler routines with the rest of the kernel
- Run them from the idle thread
 - Whenever the Linux has nothing to do it runs the idle thread
 - Waits on a semaphore until an interrupt is generated

- Essential for proper kernel functioning
 - TCP/IP timers
 - RCU synchronization
- Supported with two native operations: time and timer
- `time()` – returns the current time
- `timer()` – setups a native timer which triggers `IRQ_TIMER`
- LKL uses `NO_HZ`

- `void (*print)(const char *str, int len);`
- `long (*panic_blink)(long time);`
- `void* (*sem_alloc)(int count);`
- `void (*sem_free)(void *sem);`
- `void (*sem_up)(void *sem);`
- `void (*sem_down)(void *sem);`
- `void* (*thread_create)(void (*f)(void*), void *arg);`
- `void (*thread_exit)(void *thread);`
- `void* (*thread_id)(void);`
- `void* (*mem_alloc)(unsigned int);`
- `void (*mem_free)(void *);`
- `void (*timer)(unsigned long delta);`
- `unsigned long long (*time)(void);`
- `int (*init)(void);`
- `void (*halt)(void);`



```
static void* sem_alloc(int count)
{
    KSEMAPHORE *sem=ExAllocatePool(PagedPool, sizeof(*sem));
    if (!sem) return NULL;
    KeInitializeSemaphore(sem, count, 100);
    return sem;
}

static void sem_up(void *sem)
{
    KeReleaseSemaphore((KSEMAPHORE*)sem, 0, 1, 0);
}

static void sem_down(void *sem)
{
    KeWaitForSingleObject((KSEMAPHORE*)sem, Executive, KernelMode,
                          FALSE, NULL);
}

static void sem_free(void *sem)
{
    ExFreePool(sem);
}
```

```
static void* thread_create(void (*fn)(void*), void *arg)
{
    void *thread;
    if (PsCreateSystemThread(&thread, THREAD_ALL_ACCESS, NULL,
                            NULL, NULL, (void DDKAPI (*)(void*))fn,
                            arg) != STATUS_SUCCESS)
        return NULL;
    return thread;
}

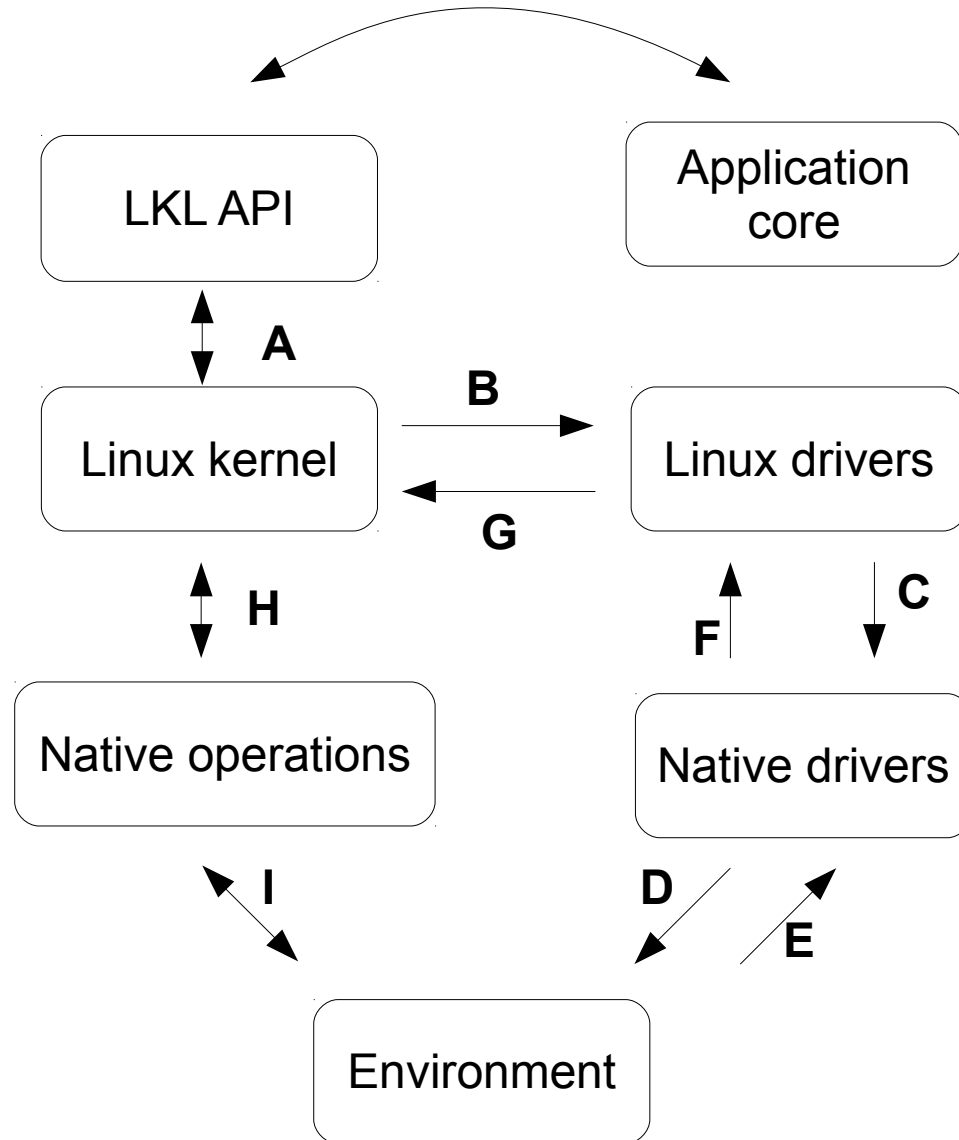
static void thread_exit(void *arg)
{
    PsTerminateSystemThread(0);
}

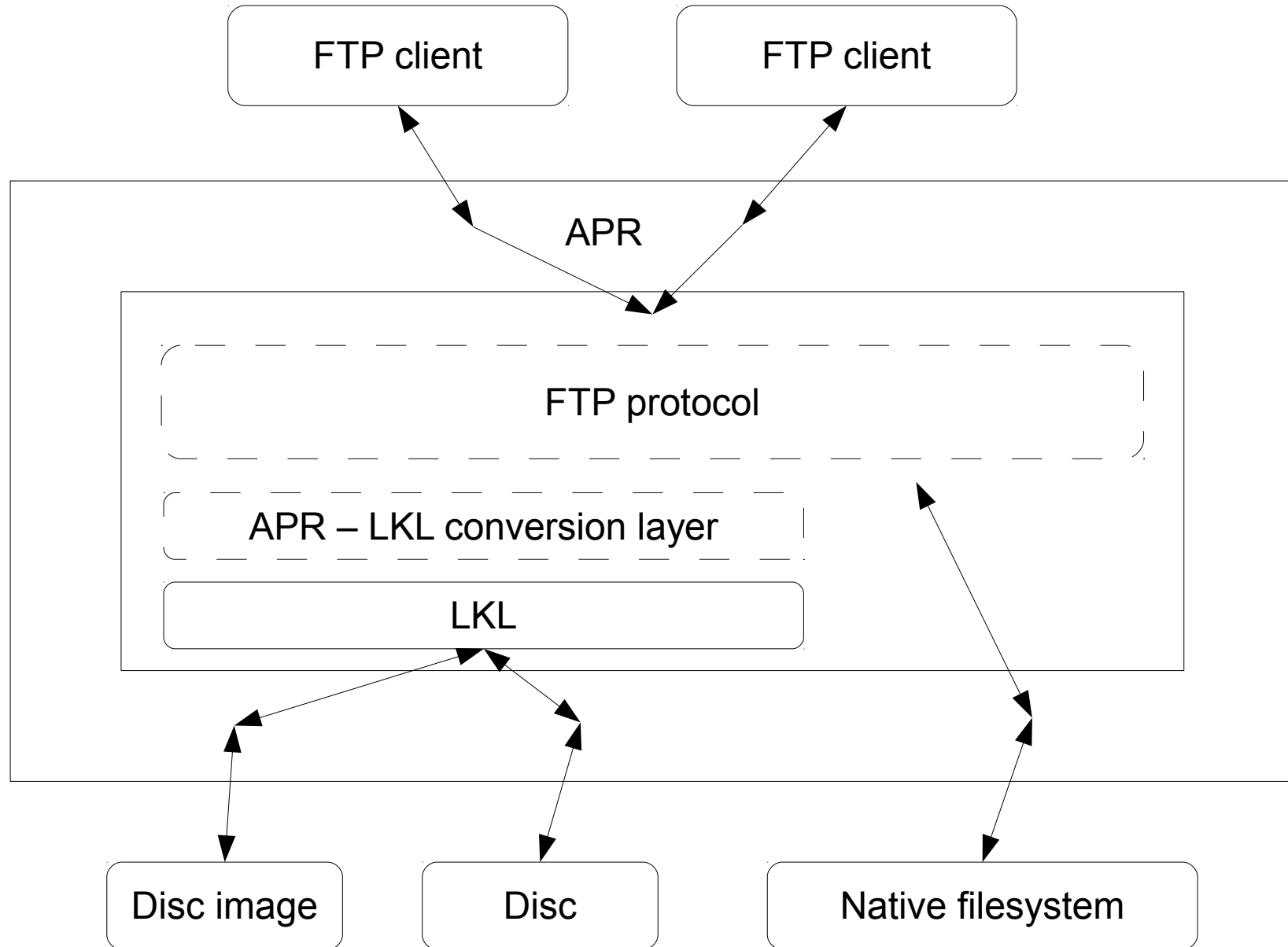
static void* mem_alloc(unsigned int size)
{
    return ExAllocatePool(NonPagedPool, size);
}

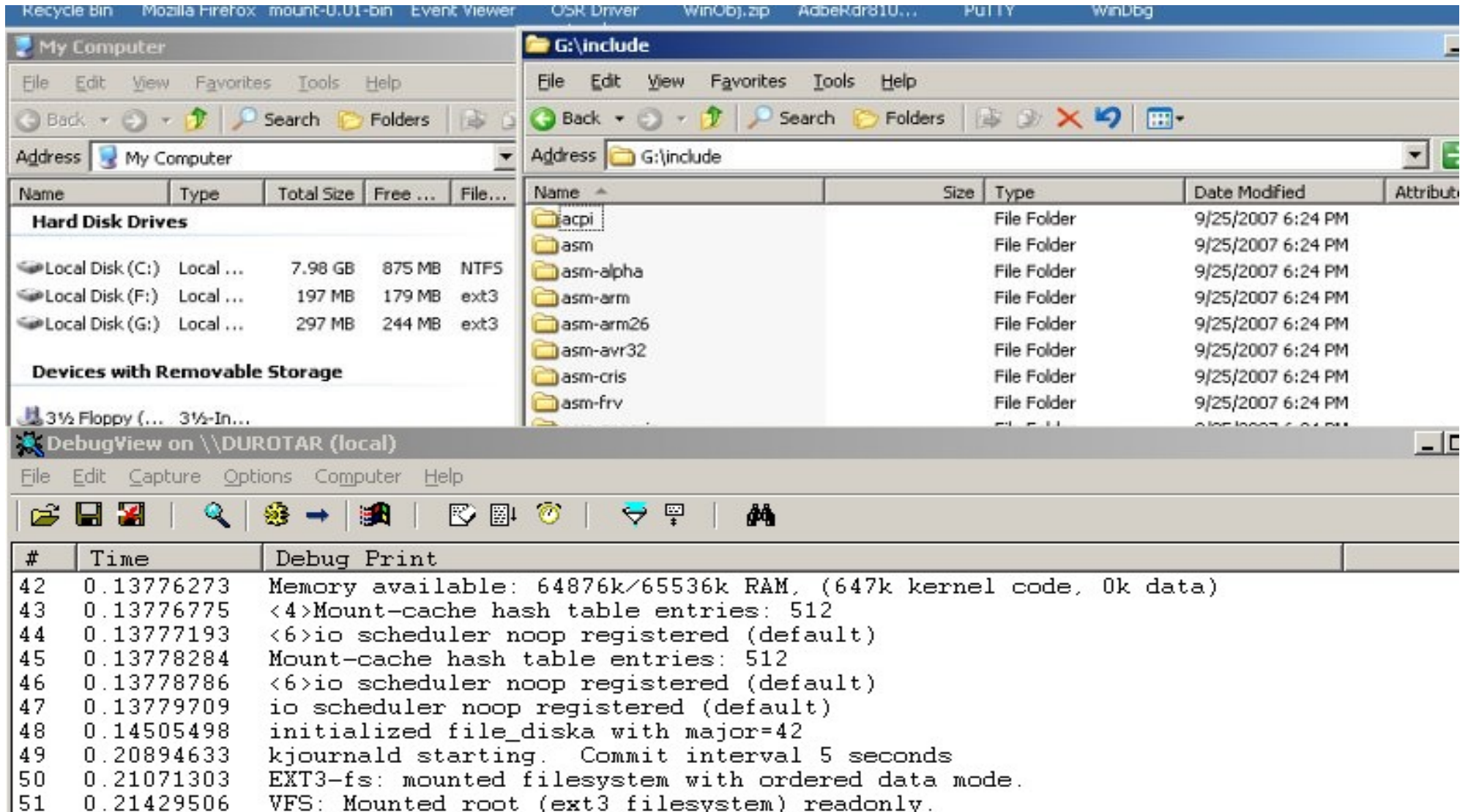
static void mem_free(void *data)
{
    ExFreePool(data);
}
```


- NT does not have an async notification mechanism
- POSIX does but we can't trigger IRQs from signal handlers
- Timer thread
 - POSIX/APR: Poll on pipe
 - NT/NTK: wait on an event object

- The application can't call directly the API – Linux functions needs to run in Linux context
- System calls
 - Application triggers `IRQ_SYSCALL`
 - The interrupt handler schedules the system call in a special kernel thread (`kyscalld`)
 - Waits on a semaphore for the system call to be finished
- In multi-threaded application only one system call can be sleeping at a time
- API to create additional syscall kernel threads







The screenshot shows a Windows desktop environment with several open windows. The primary focus is on two windows: 'My Computer' and 'G:\include'.

My Computer Window:

Name	Type	Total Size	Free ...	File...
Hard Disk Drives				
Local Disk (C:)	Local ...	7.98 GB	875 MB	NTFS
Local Disk (F:)	Local ...	197 MB	179 MB	ext3
Local Disk (G:)	Local ...	297 MB	244 MB	ext3
Devices with Removable Storage				
3 1/2 Floppy (...)	3 1/2-In...			

G:\include Window:

Name	Size	Type	Date Modified	Attribu
acpi		File Folder	9/25/2007 6:24 PM	
asm		File Folder	9/25/2007 6:24 PM	
asm-alpha		File Folder	9/25/2007 6:24 PM	
asm-arm		File Folder	9/25/2007 6:24 PM	
asm-arm26		File Folder	9/25/2007 6:24 PM	
asm-avr32		File Folder	9/25/2007 6:24 PM	
asm-cris		File Folder	9/25/2007 6:24 PM	
asm-frv		File Folder	9/25/2007 6:24 PM	

DebugView on \\DUROTAR (local) Window:

#	Time	Debug Print
42	0.13776273	Memory available: 64876k/65536k RAM, (647k kernel code, 0k data)
43	0.13776775	<4>Mount-cache hash table entries: 512
44	0.13777193	<6>io scheduler noop registered (default)
45	0.13778284	Mount-cache hash table entries: 512
46	0.13778786	<6>io scheduler noop registered (default)
47	0.13779709	io scheduler noop registered (default)
48	0.14505498	initialized file_diska with major=42
49	0.20894633	kjournald starting. Commit interval 5 seconds
50	0.21071303	EXT3-fs: mounted filesystem with ordered data mode.
51	0.21429506	VFS: Mounted root (ext3 filesystem) readonly.

- Run Valgrind's memcheck on kernel code
 - New SL*B allocator – allows Valgrind to get in the loop
 - TODO: page allocator
- „HTTP” client
 - Reuses the Linux TCP/IP stack
 - Coverage test for Linux's softirq subsystem
 - Tested on PPC
 - Native:LKL performance 4:1
- LUA-LKL
- Network simulator

- The model allows Linux code reutilization across OS, platforms, kernel/user spaces
- It allows us to keep up with the Linux change rate
- Implementing a new execution environment is easy
- Using it to develop applications is easy

<http://github.com/lkl>