

An Introduction to Multi Threaded Programming with POSIX Threads and Linux



Liam Widdowson [lbw@telstra.com]
Hewlett-Packard Consulting

January 12, 2001

Contents

1	Preface	3
1.1	General Portability	3
1.2	POSIX Threads Portability	4
2	Introduction	5
2.1	Definition	5
2.2	Brief History	5
2.3	Motivation	6
3	Thread Models	7
3.1	User Space Threads	7
3.1.1	Overview	7
3.1.2	Architecture	9
3.2	Kernel Space Threads	10
3.2.1	Overview	10
3.2.2	Architecture	11
3.3	Kernel and User Space Threads	12
3.3.1	Overview	12
3.3.2	Architecture	13
4	LinuxThreads	14
4.1	Overview	14

4.2	Simple Kernel Threads Example	15
4.3	Developing with LinuxThreads	20
5	Programming with POSIX threads	21
5.1	POSIX thread facilities	21
5.2	POSIX thread creation	22
5.2.1	Threaded version	25
5.2.2	Forking version	26
5.2.3	Performance Comparison Graph	28
5.2.4	Performance Comparison Analysis	29
5.3	Mutual Exclusion	30
5.4	Signaling and Synchronisation	34
5.5	UNIX Signal handling	40
5.6	One time operations	44
5.7	Thread Cancellation	46
5.8	Forking and Threads	51
5.9	Using Non-thread Safe Functions	52
5.10	Advanced Threading Topics Not Covered	53
6	Acknowledgments	55

1 Preface

1.1 General Portability

New open source applications are being released each day. Many of the programmers who develop open source software may do so on a single platform such as Linux. Far too much open source software is not portable out of the box. Such software typically requires non-trivial changes to compile or work reliably on other operating systems.

Portability is a key component to producing good quality, reliable software. Developing and testing software on multiple platforms helps iron out obscure bugs that many seldom be seen on a single hardware architecture or operating system.

It is a disservice to the developer's efforts if no attempt is made to make software portable. Wherever possible, a developer should use what is available in the POSIX APIs and not unnecessarily use operating system specific system calls or third party libraries. If system specific functions must be used, they should be placed in pre-processor condition statements. When designing software, thought should be given to how the software may behave on other software and hardware platforms (e.g different byte order).

Our great and fearless leader Linus Torvalds once said:

Porting this new operating systems to other platforms was really not on my mind at the beginning. At first I just wanted something that would run on my 386 [1]

It is important that developers do not fall into the same trap as Linus did. This oversight meant that a great deal of internal re-development was required to get Linux to a state where it could be ported in a straight forward manner. Most of the popular open source software has been extensively ported to other platforms (e.g Sendmail, BIND, NetBSD). It could indeed be argued that portability is a key component to the success of a particular piece of software.

1.2 POSIX Threads Portability

Portability is imperative when developing with POSIX threads. The POSIX thread API abstracts its data types from the underlying data storage structures. Developers should never make any assumptions in regards to how POSIX thread data types, the scheduler or the threading model are implemented by the vendor. All POSIX thread data types should be considered *opaque* [2].

The importance of portability will become more apparent as we delve further into POSIX threads.

2 Introduction

2.1 Definition

In a general sense a thread can be considered as a scheduling entity. More specifically consider a thread as an independent flow of control within a scheduling entity. In a typical Unix implementation, a thread will share the same address space, file descriptors, text and data segments of its parent process. Each thread has its own private stack, register context and program counter.

2.2 Brief History

The concept of a thread (as an independent flow of control) dates back to 1965 with the Berkeley Timesharing System. However, at that time they were known as processes, not threads [3]. Processes interacted with what would now be considered traditional Inter-Process Communication (IPC) means such as semaphores and message passing. One of the first thread implementations was developed on Multics at Bell Laboratories during early 1970. It used multiple stacks in a single process to support multiple background compilations [4]. In the early 1980s, micro-kernel based operating systems such as Amoeba and Chorus provided *lightweight processes* that shared the address space of a single standard process [4].

Threads as we have come to know them today were created as a part of the Open Software Foundation's Distributed Computing Environment. OSF's DCE threads is a complete user space implementation based on an early POSIX threads draft. It has been ported extensively to many operating systems thanks to its modular design. DCE could typically be ported to a new operating system in a matter of days with *some assembly required* for context switch routines [2]. Today, threads are defined by the POSIX 1003.1c-1995 standard. Extensions to the standard have been defined by The Open Group as part of the Single UNIX Specification V2. Additionally, many vendors have added non-portable system specific extensions.

2.3 Motivation

Most non-trivial applications must perform several tasks at once. Of those applications some or all of those tasks may be independent of each other and may lend themselves to parallel execution. A mail transport agent provides an excellent example - a traditional implementation uses a listener process that waits for requests and *forks* a new process to service each new client. Such an architecture has some disadvantages - forking each process adds significant overhead as *fork(2)* is an expensive system call (even if copy-on-write functionality is provided within the kernel) [5]. As each process has its own address space it must use standard IPC facilities such as semaphores or shared memory to communicate. Some of these calls are relatively expensive and retard the application's performance.

Threads address many of the aforementioned issues. Threads have a common address space which provides simplified data sharing and lightweight process creation. Additionally, threading environments provide advanced data types which assist in providing a more natural style of parallel programming.

3 Thread Models

There are three generally accepted models used in the implementation of threads on Unix variant operating systems. Each model has its own set of advantages and disadvantages. The following section examines the three models.

3.1 User Space Threads

3.1.1 Overview

A multi-threading sub system may be implemented entirely in user space. This model is typically referred to as M x 1 threads (i.e M user space threads are contained within one kernel space scheduling entity). Typically, implementations of this model are based on the POSIX threads draft 4. OSF's DCE is one such implementation. A user space library creates, terminates, schedules and synchronises threads. These threads are not directly visible to the kernel, in fact, the kernel is unaware that the single process which holds

the threads is any different from other processes.

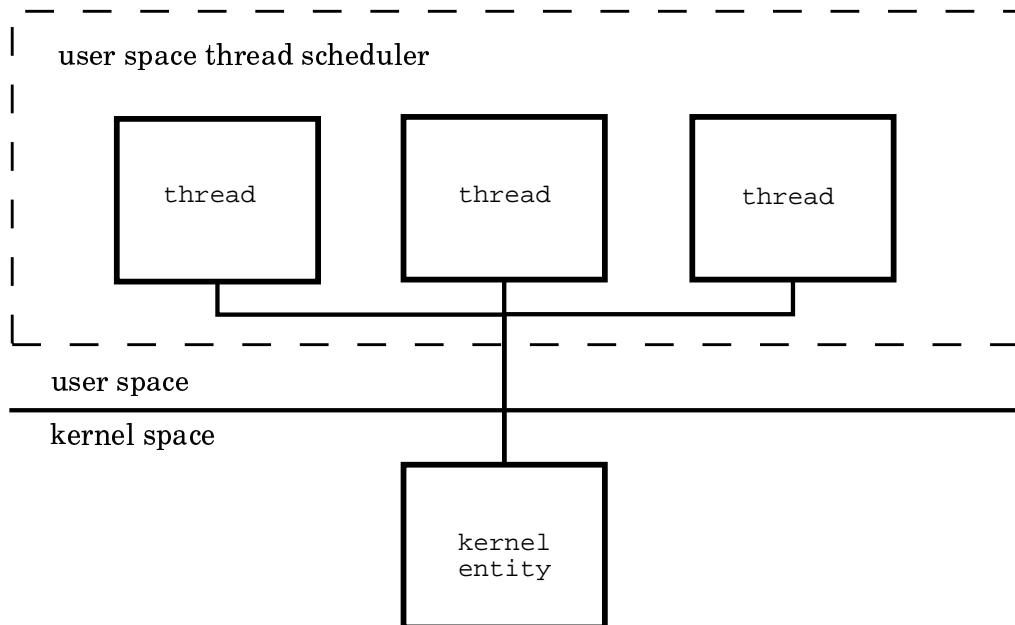
This model provides a performance benefit as each thread context switch does not traverse the user/kernel space boundary. However, there are also disadvantages associated with this model - All threads reside within a single kernel space scheduling entity (typically a Unix process). This means that a only single processor may be utilised to schedule all threads. Within uni-processor environments this may not represent a significant issue, however in multi-processor environments this model may be unacceptable.

Further, if a single thread blocks on a system call such as *read(2)*, the entire process and all associated threads block. Many user space libraries help alleviate this problem by using non-blocking I/O. However, there are still performance and compatibility issues associated with this work-around (e.g not all interfaces support non-blocking I/O).

Essentially, a user-space thread library can provide *concurrency* (i.e the appearance of multiple tasks being performed at once) but can never provide *parallelism* (i.e multiple tasks performed at the same instance in time on separate processors).

Operating systems which implement this threading model include:
FreeBSD and OpenBSD

3.1.2 Architecture



<i>Advantages</i>	<i>Disadvantages</i>
Context-switches performed entirely in user mode	Poor performance due to blocking system calls and/or associated code complexity to avoid them
Straight forward implementation	Unable to use multiple processors

3.2 Kernel Space Threads

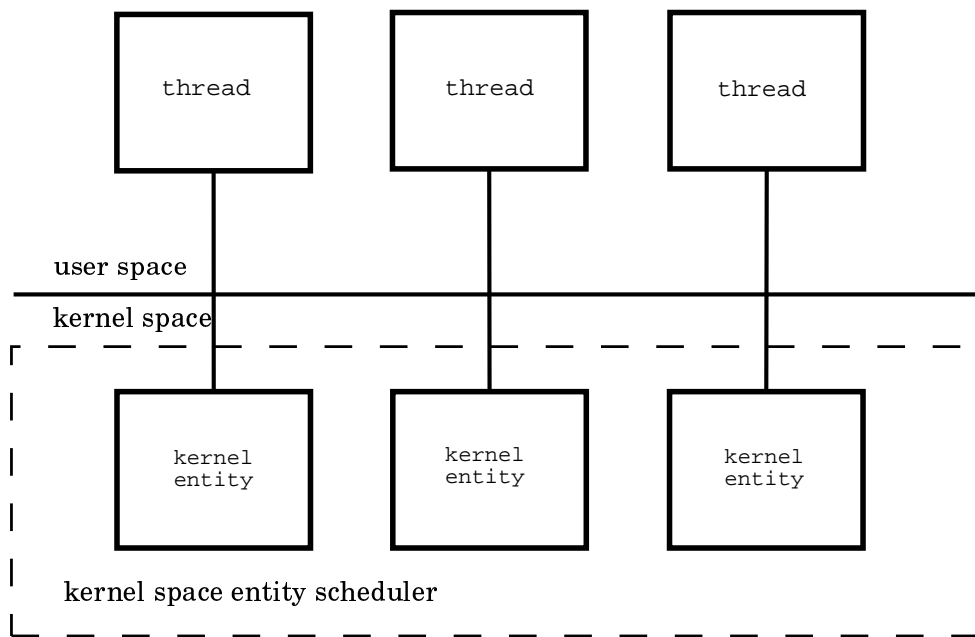
3.2.1 Overview

Unlike user space threads which do not involve any thread specific kernel interaction, kernel space threads rely on it. The kernel space thread model is often referred to as the 1 x 1 thread model (i.e one thread maps to one kernel space entity). Typically, implementations may map each single thread to a single kernel entity such as a standard Unix process or a more lightweight *kernel thread*.

Each thread is scheduled independently by the kernel scheduler so if a thread blocks it will not effect any other thread. However, creation, synchronisation and termination of kernel threads may suffer some performance problems when compared to user space threads. The kernel must be involved in operations such as requesting a memory lock or thread scheduling. Further, each kernel thread or process takes up valuable kernel resources which inhibits the creation of a large number of threads (most operating systems will only allow an arbitrary number of threads to be created system wide). Whilst the overhead may be greater, kernel threads do offer true *parallelism* as each thread may run independently on separate physical processors.

Operating systems which implement this threading model include:
HP-UX 11 and Linux

3.2.2 Architecture



<i>Advantages</i>	<i>Disadvantages</i>
Able to use multiple processors	Context switches performed entirely in kernel mode
Improved performance when coping with blocked system calls, etc	Decreased performance with large number of threads

3.3 Kernel and User Space Threads

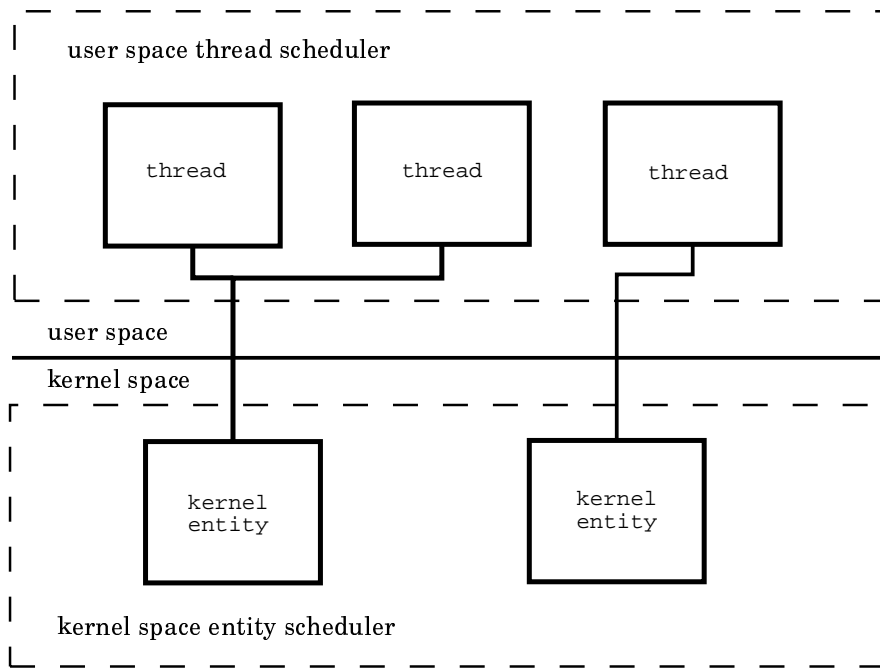
3.3.1 Overview

Kernel and user space threads are often referred to as M x N threads (i.e M user space threads within N kernel entities). The kernel and user space threading model provides a “best of both worlds” approach.

This model requires co-operation between the kernel scheduler and a user space thread library/scheduler. Essentially, sets of threads are *bound* to individual kernel entities. This allows the majority of context-switches to occur in user mode whilst still using multiple kernel entities to provide *parallelism*. For example, an application may have nine threads which are spread across four kernel entities. This would allow the application to take advantage of multi-processor hardware whilst still traversing the user-kernel boundary as little as possible. In situations where a thread blocks in a system call, another thread is immediately scheduled thus not halting the execution of threads that are bound to the same kernel entity.

Operating systems which implement this threading model include:
IRIX, Solaris and Tru64

3.3.2 Architecture



<i>Advantages</i>	<i>Disadvantages</i>
Able to use multiple processors	Complex system implementation
Improved performance with blocked system calls, etc	
Most context switches in user mode	

4 LinuxThreads

4.1 Overview

Linux implements the 1 x 1 kernel space threading model. Thread creation at a kernel level is performed with *clone(2)* system call. Under Linux, the *clone(2)* system call is a generalisation of *fork(2)*. Thus a thread under Linux is not significantly different from a standard heavy-weight Unix process. When a process creates its first thread, a total of three threads will be active. The initial or *main()* thread (process), the newly created thread and a manager thread which helps performs scheduling activities. This is significantly different from implementations such as Solaris and Tru64 which have separate kernel and user space entities to represent threads.

Linux's kernel developers have made a specific decision to support this model as they believe it removes the complexity associated with implementing another entity and scheduler within the kernel. Threads are created, terminated and scheduled by the same code that does so for processes. Apart from providing simplicity, the rationale behind this architecture choice is that Linux has exceptionally fast context-switch code and this negates what is one of the biggest problems associated with 1 x 1 implementations - context switch performance. However, this simplicity comes with an inherent problem - Linux will never truly be POSIX 1003.1c-1995 compliant without specific support within the kernel for threads.

It is for this reason that threads under Linux do not provide a significant performance benefit as they do on other platforms. Nonetheless, threads provide a more natural parallel style suitable for parallel programming. However, keep in mind portability - a threaded application will perform significantly better than a forking application on most other operating systems.

4.2 Simple Kernel Threads Example

The following source code was posted by Linus Torvalds to the Linux kernel mailing list in 1996. It provides a basic example of how the *clone(2)* system call may be used to implement kernel threads.

```
_____ linux-clone.c - Linus Torvalds _____  
1 #include <signal.h>  
2 #include <stdio.h>  
3 #include <stdlib.h>  
4 #include <fcntl.h>  
5 #include <linux/unistd.h>  
6  
7 #define STACKSIZE 16384  
8  
9 /* signal mask to be sent at exit */  
10 #define CSIGNAL      0x000000ff  
11 /* set if VM shared between processes */  
12 #define CLONE_VM     0x00000100  
13 /* set if fs info shared between processes */
```



```

14 #define CLONE_FS          0x00000200
15 /* set if open files shared between processes */
16 #define CLONE_FILES      0x00000400
17 /* set if signal handlers shared */
18 #define CLONE_SIGHAND    0x00000800
19
20 int start_thread(void (*fn)(void *), void *data) {
21     long retval;
22     void **newstack;
23
24     /*
25      * allocate new stack for subthread
26      */
27     newstack = (void **) malloc(STACKSIZE);
28     if (!newstack)
29         return -1;
30     /*
31      * Set up the stack for child function, put the (void *)
32      * argument on the stack.
33      */
34     newstack = (void **) (STACKSIZE + (char *) newstack);
35     *--newstack = data;
36
37     /*
38      * Do clone() system call. We need to do the low-level stuff
39      * entirely in assembly as we're returning with a different
40      * stack in the child process and we couldn't otherwise guarantee
41      * that the program doesn't use the old stack incorrectly.

```

```

42      *
43      * Parameters to clone() system call:
44      *      %eax - __NR_clone, clone system call number
45      *      %ebx - clone_flags, bitmap of cloned data
46      *      %ecx - new stack pointer for cloned child
47      *
48      * In this example %ebx is CLONE_VM | CLONE_FS | CLONE_FILES |
49      * CLONE_SIGHAND which shares as much as possible between parent
50      * and child. (We or in the signal to be sent on child
51      * termination into clone_flags: SIGCHLD makes the cloned
52      * process work like a "normal" unix child process)
53      *
54      * The clone() system call returns (in %eax) the pid of the newly
55      * cloned process to the parent, and 0 to the cloned process. If
56      * an error occurs, the return value will be the negative errno.
57      *
58      * In the child process, we will do a "jsr" to the requested
59      * function and then do a "exit()" system call which will
60      * terminate the child.
61      */
62      __asm__ __volatile__(
63          "int $0x80\n\t"           // Linux/i386 system call
64          "testl %0,%0\n\t"        // check return value
65          "jne 1f\n\t"             // jump if parent
66          "call *%3\n\t"           // start subthread function
67          "movl %2,%0\n\t"
68          "int $0x80\n\t"         // exit syscall: ex subthread
69          "1:\t"

```

```

70         :="a" (retval)
71         :="0" (__NR_clone),"i" (__NR_exit),
72         "r" (fn),
73         "b" (CLONE_VM | CLONE_FS | CLONE_FILES | CLONE_SIGHAND |
74             SIGCHLD),
75         "c" (newstack));
76
77     if (retval < 0) {
78         errno = -retval;
79         retval = -1;
80     }
81     return retval;
82 }
83
84 int show_same_vm;
85
86 void cloned_process_starts_here(void * data) {
87     printf("child:\t got argument %d as fd\n", (int) data);
88     show_same_vm = 5;
89     printf("child:\t vm = %d\n", show_same_vm);
90     close((int) data);
91 }
92
93 int main() {
94     int fd, pid;
95
96     fd = open("/dev/null", O_RDWR);
97     if (fd < 0) {

```

```

98         perror("/dev/null");
99         exit(1);
100     }
101     printf("mother:\t fd = %d\n", fd);
102
103     show_same_vm = 10;
104     printf("mother:\t vm = %d\n", show_same_vm);
105     pid = start_thread(cloned_process_starts_here, (void *) fd);
106
107     if (pid < 0) {
108         perror("start_thread");
109         exit(1);
110     }
111
112     sleep(1);
113     printf("mother:\t vm = %d\n", show_same_vm);
114     if (write(fd, "c", 1) < 0)
115         printf("mother:\t child closed our file descriptor\n");
116 }

```

linux-clone.c

The output from Linus' threading example is as follows:

```

mother: fd = 3
mother: vm = 10
child:  got argument 3 as fd
child:  vm = 5
mother: vm = 5
mother: child closed our file descriptor

```

4.3 Developing with LinuxThreads

A variety of threading libraries (including user space implementations) have been developed for Linux. However, the defacto standard library is now LinuxThreads. LinuxThreads 0.8 provides a mostly compliant POSIX 1003.1c implementation and forms part of glibc 2.x. Linux initially diverges from the POSIX standard with one fundamental difference - each thread has an independent PID. This occurs because each thread is essentially a standard Unix process pretending to be a thread. POSIX requires that all threads in a process share a single PID. In turn, this violation effects the behavior of *fork(2)* and signals, none of which are entirely POSIX compliant.

The following are minimum requirements for developing threaded software with Linux: gcc 2.8, gdb 4.18, LinuxThreads 0.8, glibc 2.0.1 and kernel 2.x

The following table provides examples of common linker and compiler flags required when compiling POSIX threaded software:

<i>Utility</i>	<i>Parameter</i>	<i>Platform</i>
ld	-lpthread	All platforms
gcc/cc	-D_REENTRANT -pthread -D_POSIX_C_SOURCE=199506L -mt (SunSoft CC only)	All platforms Free/OpenBSD HP-UX Solaris

5 Programming with POSIX threads

This Section outlines how POSIX threads functions and data types may be used to create multi threaded applications. Within Section 5.2 performance comparisons will be made between a forking and threaded example application on Linux, Solaris and FreeBSD. This comparison is by no means a definitive benchmark, however, it does provide some insight into the performance benefits of threaded applications. The following systems were used to test each piece of code.

<i>Hostname</i>	<i>Operating System</i>	<i>Processor(s)</i>	<i>Memory</i>
athena	Linux 2.2.17	AMD K6-2/300MHz	64MB
fatso	FreeBSD 4.2	Intel 100MHz Pentium	32MB
helios	Solaris 8	160MHz TurboSPARC	96MB
prometheus	Solaris 8	2 x 150MHz HyperSPARC	256MB

5.1 POSIX thread facilities

In addition to providing facilities to create, terminate and synchronise threads and shared data, POSIX threads implementations should provide features such as:

- A per-thread *errno* global variable;
- Thread safe versions of functions such as *malloc(3)* and *free(3)*.

5.2 POSIX thread creation

A POSIX thread is created using the *pthread_create()* function. This function is defined as:

```
int pthread_create(pthread_t *thread_id, pthread_attr_t *attr,
                  void *(*startroutine)(void *), void *arg)
```

The *pthread_create()* function returns 0 upon success, ENOMEM if there are insufficient resources to create the thread, EINVAL if the attributes are invalid or EPERM if there are insufficient privileges to create the thread.

The *thread_id* pointer stores the thread's unique identifier (analogous to a process PID), *attr* stores the attributes that the new thread should possess upon creation and *startroutine* points to the function that should be called upon thread creation. Finally, *arg* is a pointer to data which should be supplied as an argument to *startroutine*. Note, only one argument is allowed so multiple data objects may be placed within a struct.

Functions which are called by *pthread_create()* should be of the form:

```
void *function(void *argument);
```

The argument may then be re-cast at a later date to the desired type. When a thread wishes to exit, it may do so in a variety of ways - either via a *return* statement or by calling the following function:

```
void pthread_exit(void *value_ptr);
```

The value of the exiting thread is lost unless the caller waits for it with the following function:

```
void pthread_join(pthread_t thread_id, void *value_ptr);
```

The *pthread_join()* function is analogous to *waitpid(2)* for processes.

Thread attributes are set by the *pthread_attr_**(*)* routines and define details such as what state the thread should begin execution in, scheduling information, etc. If NULL is specified as an argument to *pthread_create()*, the thread will be created with a default set of attributes. An examination of all pthread attributes is beyond the scope of this paper, but one noteworthy attribute must be mentioned. When a thread exits, its resources are not automatically *recycled*. The thread will continue to exist in a zombie state until it is joined with *pthread_join()*.

However, in some situations it may not be convenient to use *pthread_join()*. The thread's id may be lost, the caller may not wish to collect the thread's return value or suspend its execution while it waits for the thread to finish. In such a case, the thread may be created with the PTHREAD_CREATE_DETACHED attribute set. This will create the thread in a detached state. When the thread exits all resources will be automatically *recycled*.

The attribute may be set with code such as the following or by calling the `pthread_detach(pthread_t *thread_id)` function at any time:

```
pthread_attr_t attr;
pthread_t      t;

pthread_attr_init(&attr);
pthread_attr_setdetachstate(&attr, PTHREAD_CREATE_DETACHED);
pthread_create(&t, &attr, some_function(), (void *) some_arg);
..
or
..
pthread_detach(&t1)
```

A thread that is detached can never be joined again with `pthread_join()`. It is also worth mentioning again, that all pthread data types are *opaque*. One should never make any assumptions about the underlying data types. Further, the result of copying pthread data variables is undefined as is performing arithmetic (e.g equality) operations upon them. In order to check if two `pthread_t` thread IDs are the same, the following function should be used:

```
int pthread_equal(pthread_t thread0, pthread_t thread1);
```

The `pthread_equal()` function returns a non-zero value if the thread IDs are equal, otherwise zero is returned.

A thread's thread identifier may be retrieved at any time by using the following function:

```
pthread_t pthread_self(void);
```

The following code demonstrates thread creation by creating a pre-defined number of threads and prints a message from each. To contrast, a forking version is provided.

5.2.1 Threaded version

```
_____ test2-thread.c _____  
1 #include <pthread.h>  
2 #include <stdio.h>  
3 #include <string.h>  
4 #include <errno.h>  
5 #include <unistd.h>  
6 #define MAX_THREADS 512  
7  
8 void *say_hello(void) {  
9     printf("[thread %d] new thread created\n", (int) pthread_self());  
10    sleep(1);  
11    return NULL;  
12 }  
13  
14 int main(void) {  
15    pthread_t    t1;
```

```

16         int                x, status;
17
18         #ifdef __sun__
19         pthread_setconcurrency( sysconf(_SC_NPROCESSORS_ONLN)+1 );
20         #endif
21
22         for (x = 0; x < MAX_THREADS; x++) {
23             status = pthread_create(&t1, NULL, say_hello, NULL);
24             if (status != 0) {
25                 fprintf(stderr, "[thread %d] error creating thread %s\n",
26                     (int) pthread_self(), (char *) strerror(status));
27                 exit(0);
28             }
29         }
30         printf("[thread %d] this is the end\n", (int) pthread_self());
31         exit(0);
32     }

```

test2-thread.c

5.2.2 Forking version

```

1 #include <stdio.h>
2 #include <unistd.h>
3 #include <errno.h>
4 #include <sys/types.h>
5 #include <sys/wait.h>
6 #define MAX_PROCESSES 512

```

```

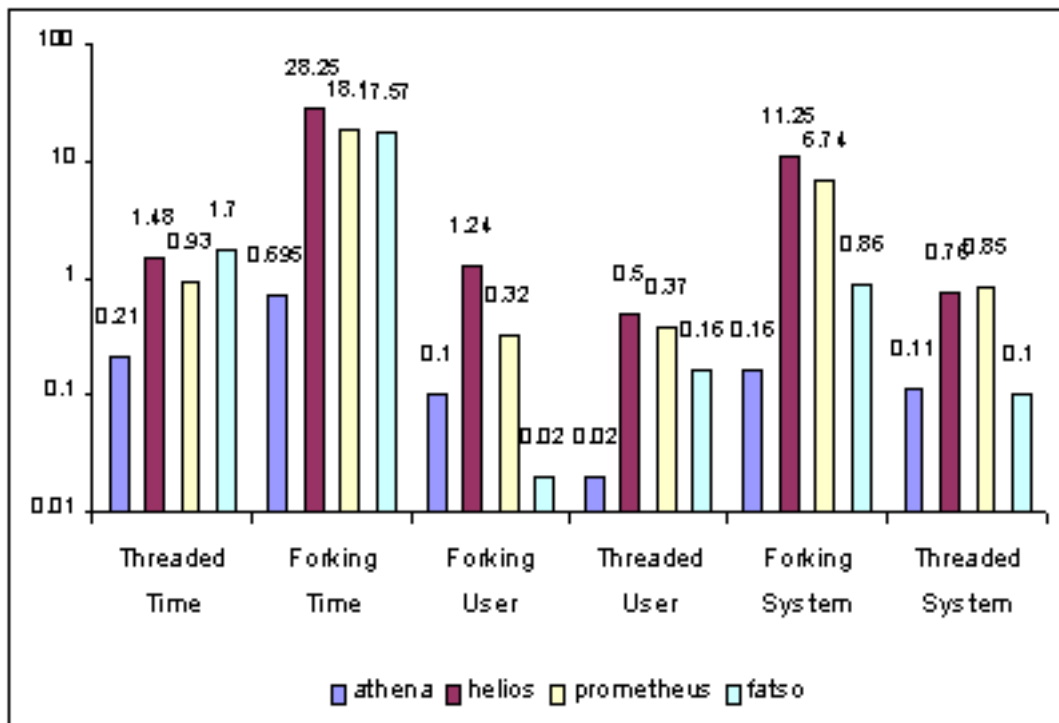
7
8 void say_hello(void) {
9     printf("[pid %d] I am a new child\n", (int) getpid());
10    sleep(1);
11 }
12
13 int main(void) {
14     int    x;
15     pid_t  pid;
16     for (x = 0; x < MAX_PROCESSES; x++) {
17         pid = fork();
18         if (pid == (pid_t) -1) {
19             perror("fork");
20             exit(0);
21         }
22         if (pid == (pid_t) 0) {
23             // in the child
24             say_hello();
25             exit(0);
26         } else {
27             // in the parent
28         }
29     }
30     printf("[pid %d] this is the end\n", (int) getpid());
31     return 0;
32 }

```

test2-process.c

5.2.3 Performance Comparison Graph

The following graph displays the time taken by the threaded and forking test code. Note that the Y axis represents time using a logarithmic scale.



5.2.4 Performance Comparison Analysis

It is evident from the data provided that threaded applications perform significantly better. On Solaris, the threaded application performed its tasks almost twenty times faster than the forking version. It spent almost half as much time in user mode and sixteen times less time in kernel mode. Attention should be paid to the difference in results between the single and multi-processor Solaris machines. The multi-processor machine spent slightly more time in kernel mode and slightly less in user mode. This is because some of the threads could run on different processors thus slightly more kernel intervention was required for context switching and thread creation.

On Linux, the gap was much smaller, but threads still provide a performance advantage finishing their tasks three times faster whilst spending five times less time in user mode and about the same time in kernel mode.

FreeBSD's user space threads implementation also provided excellent results. The threaded application performed almost ten times faster and spent significantly less time in both user and kernel mode.

The code used in this simple experiment does not mimic the behavior of a typical application (which would be both computational and I/O bound). In such a situation a threading model which includes a kernel and user space scheduler (such as Solaris or Tru64) would most likely perform best.

5.3 Mutual Exclusion

As threads share the same address space, they must possess a synchronisation method to ensure that shared resources are not corrupted by concurrent reading and writing. POSIX threads provide a synchronisation type called a mutex. A mutex (short for mutual exclusion) may be used to protect access to any resource. A mutex locked within a single thread may not be unlocked by any other thread. A thread which attempts to acquire a locked mutex will block until it is released.

The mutex type within POSIX threads is *pthread_mutex_t* and may be initialized at deceleration with the PTHREAD_MUTEX_INITIALIZER value or with the following function:

```
int pthread_mutex_init( pthread_mutex_t *mutex,
                       pthread_mutexattr_t *attr);
```

The mutex initialization returns 0 if successful, EBUSY if it has already been initialised or EINVAL if the mutex or attributes are invalid.

A mutex may have a variety of attributes accrued to it at initialisation. Discussion of these attributes is beyond the scope of this paper.

Mutexes should be used to protect multiple threads from reading *or* writing shared data at once. In addition, a mutex should be used to protect a non-thread safe function from multiple invocations across threads.

The following code provides an example of how a mutex is used to protect stdout. When a program is linked with the pthread library, the library will transparently perform internal mutex locking for functions such as *malloc(3)* and buffering/locking of stdout. In the following code explicit mutex locking of stdout is required as characters are being written to stdout individually.

```
----- file-mutex.c -----
1 #include <pthread.h>
2 #include <stdio.h>
3 #include <time.h>
4 #include <string.h>
5 #include <errno.h>
6 #include <unistd.h>
7
8 pthread_mutex_t stdout_mutex = PTHREAD_MUTEX_INITIALIZER;
9 const char array[] = "http://www.linux.conf.au/\n";
10
11 void *display(void) {
12     int x, y;
13     struct timespec tv = { 0, 100 };
14     for(y=0;y<5;y++) {
15         pthread_mutex_lock(&stdout_mutex);
16         printf("[thread %d] ", (int) pthread_self());
17         for (x=0;x<strlen(array);x++)
18             putchar(array[x]);
19         pthread_mutex_unlock(&stdout_mutex);
```



```

20         nanosleep(&tv, NULL);
21     }
22     return NULL;
23 }
24
25 int main(void) {
26     pthread_t    t1, t2;
27     int          status;
28     #ifdef __sun__
29     pthread_setconcurrency( sysconf(_SC_NPROCESSORS_ONLN)+1 );
30     #endif
31     status = pthread_create(&t1, NULL, (void *) display, NULL);
32     if (status != 0) {
33         fprintf(stderr, "[thread %d] error creating thread: %s\n",
34                 pthread_self(), (char *) strerror(status));
35         exit(0);
36     }
37     status = pthread_create(&t2, NULL, (void *) display, NULL);
38     if (status != 0) {
39         fprintf(stderr, "[thread %d] error creating thread: %s\n",
40                 pthread_self(), (char *) strerror(status));
41         exit(0);
42     }
43     pthread_join(t2, NULL);
44     return 0;
45 }

```

file-mutex.c

The following was output from the aforementioned program:

```
[thread 4] http://www.linux.conf.au/  
[thread 5] http://www.linux.conf.au/  
[thread 4] http://www.linux.conf.au/  
[thread 5] http://www.linux.conf.au/  
[thread 4] http://www.linux.conf.au/  
[thread 4] http://www.linux.conf.au/  
[thread 5] http://www.linux.conf.au/  
[thread 4] http://www.linux.conf.au/  
[thread 5] http://www.linux.conf.au/  
[thread 5] http://www.linux.conf.au/
```

The following was output from the aforementioned program with mutex locking removed at lines 16 and 20. Observe that how without correct synchronisation the data is corrupted.

```
[thread 5] [thread 4] http:ht//wwtp://www.w.linux.conf.linu.au/x.c  
onf.au/  
[thread 5] http://www.linux.conf.au/  
[thread 4] http://www[thread 5] http://www.linux.conf..linuxau/  
.conf.au/  
[thread 5] http://w[thread 4] hww.littnux.conf.au/  
p://www.linux.conf.au/  
[thread 4] http://www.linux.conf.au/  
[thread 5] http://www.linux.conf.au/  
[thread 4] http://www.linux.conf.au/
```

5.4 Signaling and Synchronisation

Within an application, there may be times where the need arises to signal another process or thread to begin doing some work. There may be a queue of data that needs to be processed or a thread pool waiting to service a newly established connection. POSIX threads offers *condition variables* to perform this function. A thread or many threads may wait on a particular condition and once the condition is met, one or all of the threads will return from their blocked state and begin their work.

Within POSIX threads, a condition variable is always associated with a mutex and indirectly, associated with the data it is waiting upon. A condition variable is defined by the *pthread_cond_t* type. Like a mutex, a *pthread_cond_t* can be initialized either with the constant `PTHREAD_COND_INITIALIZER` or the following function:

```
pthread_cond_init(pthread_cond_t *cond, pthread_condattr_t *attr);
```

The function *pthread_cond_init()* returns 0 if successful, `EAGAIN` if insufficient resources are available or `ENOMEM` if insufficient memory is available.

Like most of the POSIX thread functions, condition variables have attributes which may be set upon initialisation. Discussion of these attributes are beyond the scope of this paper.

Once initialised, a condition variable may be waited upon by using either of the following functions:

```
int pthread_cond_wait(pthread_cond_t *cond, pthread_mutex_t *mutex);

int pthread_cond_timedwait(pthread_cond_t *cond,
                           pthread_mutex_t *mutex,
                           struct timespec *wait_time);
```

The above functions will return 0 upon success, `EINVAL` if either the mutex or condition variable are invalid or `ETIMEDOUT` if *pthread_cond_timedwait()* returns from the timed wait without the condition variable changing state.

Before calling either of the above functions, a thread must have the associated data mutex locked. The function will internally unlock the mutex to allow other operations to occur, but once the condition has been signaled the function will return with the mutex locked.

In order to wake a thread or threads waiting upon a condition variable, one of the following functions may be used:

```
int pthread_cond_signal(pthread_cond_t *cond);

int pthread_cond_broadcast(pthread_cond_t *cond);
```

The *pthread_cond_signal()* function will only wake a single thread whilst *pthread_cond_broadcast()* will wake all waiting threads. The above functions

return 0 upon success or EINVAL if the condition variable has not been initialised. The following code demonstrates how a condition variable may be used:

```
cond-wait.c
1 #include <pthread.h>
2 #include <stdio.h>
3 #include <time.h>
4 #include <string.h>
5 #include <errno.h>
6 #include <stdlib.h>
7 #include <unistd.h>
8
9 typedef struct {
10     pthread_mutex_t mutex;
11     pthread_cond_t cond;
12     char *data;
13 } data_t;
14
15 data_t data;
16
17 pthread_mutex_t stdout_mutex = PTHREAD_MUTEX_INITIALIZER;
18
19 void *wait_thread(void) {
20     pthread_mutex_lock(&data.mutex);
21     while (1) {
22
23         /* mutex is always returned locked, hence no need to
24            lock or unlock within the loop */
```

```

25
26         pthread_cond_wait(&data.cond, &data.mutex);
27         pthread_mutex_lock(&stdout_mutex);
28
29         if (strcmp(data.data, "q\n") == 0) {
30             printf("[thread %d] goodbye!\n",
31                    (int) pthread_self());
32         } else {
33             printf("[thread %d] user said: %s",
34                    (int) pthread_self(), data.data);
35         }
36         pthread_mutex_unlock(&stdout_mutex);
37         free(data.data);
38         data.data = NULL;
39     }
40 }
41
42 void *input_thread(void) {
43     char buf[80];
44     struct timespec tv = { 0, 100 };
45     while (strcmp(buf, "q\n") != 0) {
46         pthread_mutex_lock(&stdout_mutex);
47         printf("[thread %d] enter string: ", (int) pthread_self());
48         fgets(buf, 79, stdin);
49         pthread_mutex_unlock(&stdout_mutex);
50         pthread_mutex_lock(&data.mutex);
51
52         if (data.data == NULL) {

```

```

53         data.data = strdup(buf);
54     } else {
55         free(data.data);
56         data.data = strdup(buf);
57     }
58
59     pthread_mutex_unlock(&data.mutex);
60     pthread_cond_signal(&data.cond);
61
62     nanosleep(&tv, NULL);
63 }
64 return NULL;
65 }
66
67 int main(void) {
68
69     pthread_t    t1, t2;
70     int          status;
71
72     // initialize the data
73     pthread_mutex_init(&data.mutex, NULL);
74     pthread_cond_init(&data.cond, NULL);
75     data.data = NULL;
76
77     #ifdef __sun__
78     pthread_setconcurrency( sysconf(_SC_NPROCESSORS_ONLN)+1 );
79     #endif
80

```

```

81     status = pthread_create(&t1, NULL, (void *) wait_thread, NULL);
82     if (status != 0) {
83         fprintf(stderr, "[thread %d] error creating thread: %s\n",
84                 pthread_self(), (char *) strerror(status));
85         exit(0);
86     }
87     status = pthread_create(&t2, NULL, (void *) input_thread, NULL);
88     if (status != 0) {
89         fprintf(stderr, "[thread %d] error creating thread: %s\n",
90                 pthread_self(), (char *) strerror(status));
91         exit(0);
92     }
93
94     pthread_join(t2, NULL);
95     return 0;
96 }

```

cond-wait.c

The following was output from the aforementioned program:

```

[thread 5] enter string: hello
[thread 4] user said: hello
[thread 5] enter string: world
[thread 4] user said: world
[thread 5] enter string: q
[thread 4] goodbye!

```


5.5 UNIX Signal handling

Signal handling can become complex within a multi-threaded application and should, where possible, be avoided. The POSIX thread standard states that application generated signals are sent to any thread within that application. This means that if a SIGCHLD is raised by a child process, it may not be delivered to the thread that originally created that child. However, the signals SIGFPE, SIGSEGV, SIGPIPE and SIGTRAP are always delivered to the thread that caused them [2].

Despite signals being sent to individual threads, they do effect the process as a whole. Unless the relevant signal handler is installed, a signal such as SIGSEGV (whether sent to an individual thread or the process) results in the entire process being killed and a core file generated.

As mentioned previously, LinuxThreads signal handling differs significantly from the POSIX standard. According to the POSIX standard, external signals generated from a command such as *kill* should be sent to the process which in turn delivers it to any thread that does not block the signal. Since a thread within Linux is in fact a process, the external signal will be delivered to that particular process. If another thread is blocked in *sigwait(2)* waiting for that external signal, it will remain blocked. Additionally, SIGUSR1 and SIGUSR2 may not be used within a threaded application under Linux as they are used internally by LinuxThreads.

On POSIX compliant systems the preferred signal handling method is to create a thread specifically for signal handling. The required signals should be masked in all other threads by using the following function:

```
int pthread_sigmask(int how, sigset_t *set, sigset_t *oset);
```

The *pthread_sigmask()* function returns 0 upon success and a non-zero value upon error. A thread inherits the signal masks of its creator. It is therefore possible to set signal masks of all threads by specifying them in *main()* before creating any threads.

Once all threads other than the signal handling thread have the signal masked, the handler thread may wait for the signal with *sigwait()* and perform the appropriate processing when necessary.

However, in order to process external signals in a portable fashion on Linux, a signal handler should be installed process wide with either *signal(2)* or *sigaction(2)*. The code executed within a signal handler is not async-safe so it is generally necessary to wake a thread (using Sys V IPC) to perform the required operations once a signal has been received.

The following code outlines how a portable external signal handler may be implemented:

```
----- signals.c -----
1 #include <pthread.h>
2 #include <stdio.h>
3 #include <string.h>
4 #include <signal.h>
5 #include <sys/types.h>
6 #include <unistd.h>
7 #include <semaphore.h>
8
9 sem_t          thread_sem;
10
11 void *signal_thread(void) {
12     while (1) {
13         // wait for the semaphore to be posted
14         sem_wait(&thread_sem);
15         printf("[thread %d] SIGHUP received\n",
16             (int) pthread_self());
17     }
18 }
19
20 void handler(int sig) {
21     if (sig == SIGHUP) {
22         signal(SIGHUP, handler);
23         sem_post(&thread_sem);
24     }
25 }
26
27 int main(void) {
28
```

```

29     pthread_t      t;
30     int            status;
31
32     #ifdef __sun__
33     pthread_setconcurrency( sysconf(_SC_NPROCESSORS_ONLN)+1 );
34     #endif
35
36     if ( sem_init(&thread_sem, 0, 0) == -1) {
37         printf("[thread %d] error initialising semaphore: ",
38             (int) pthread_self());
39         perror("sem_init");
40         exit(1);
41     }
42
43     signal(SIGHUP, handler);
44
45     status = pthread_create(&t, NULL, (void *) signal_thread, NULL);
46
47     if (status != 0) {
48         printf("[thread %d] error creating signal_thread: %s\n",
49             (int) pthread_self(), strerror(status));
50         exit(1);
51     }
52
53     pthread_join(t, NULL);
54     return 0;
55 }

```

signals.c

5.6 One time operations

Sometimes within a threaded program, it is useful to execute an operation only once within a thread, no matter how many times it is invoked. This may include operations such as initialising a mutex with the *pthread_mutex_init()* function. This may be performed by declaring a variable of type *pthread_once_t* which must be initialized with the value `PTHREAD_ONCE_INIT`. It is then possible to perform an operation a guaranteed single time by using the following function:

```
int pthread_once(pthread_once_t *once, void (*once_routine, void));
```

once_routine is the pointer of the function to be called. Note, no argument can be given to the function unlike with *pthread_create_thread()*. The function *pthread_once()* returns 0 if successful or `EINVAL` if either the supplied routine or *pthread_once_t* are invalid.

The following code illustrates how *pthread_once()* may be used:

```
----- do-once.c -----  
1 #include <pthread.h>  
2 #include <stdio.h>  
3 #include <string.h>  
4 #include <errno.h>  
5 #include <unistd.h>  
6 #define MAX_THREADS 512  
7 pthread_once_t once = { PTHREAD_ONCE_INIT };
```

```

8 void say_hello(void) {
9     printf("[thread %d] Listen carefully, I will say this only once\n",
10          (int) pthread_self());
11 }
12
13 void *thr_ctl(void) {
14     pthread_once(&once, say_hello);
15     return NULL;
16 }
17
18 int main(void) {
19     pthread_t      t1;
20     int            x, status;
21     #ifdef __sun__
22     pthread_setconcurrency( sysconf(_SC_NPROCESSORS_ONLN)+1 );
23     #endif
24     for (x = 0; x < MAX_THREADS; x++) {
25         status = pthread_create(&t1, NULL, thr_ctl, NULL);
26         if (status != 0) {
27             fprintf(stderr, "[thread %d] error creating thread %s\n",
28                  (int) pthread_self(), (char *) strerror(status));
29             exit(0);
30         }
31     }
32     printf("[thread %d] This is the end\n", (int) pthread_self());
33     exit(0);
34 }

```

do-once.c

5.7 Thread Cancellation

In many situations, a programmer may wish to stop a thread's operation and ensure that it shuts down cleanly. An example would be ensuring a dbm database is closed correctly with *dbm_close(3)* before an application is shut down or unlocking shared mutexes when the thread is shutdown or restarted. It would be possible, but messy to implement such a signaling method using UNIX signals or condition variables. Fortunately, POSIX threads provides a set of cancellation functions which may be used to shutdown a thread cleanly.

In order for a thread to be canceled, it must first declare that it is able to be canceled by calling the following functions:

```
int pthread_setcancelstate(int state, int *oldstate);
```

```
int pthread_setcanceltype(int type, int *oldtype);
```

POSIX threads provides two types of cancellation - deferred and asynchronous. A thread using deferred cancellation will only allow cancellation at specific cancellation points. Cancellation points include various system and function calls including *read(2)*, *pthread_cond_wait()*, etc. Consult the relevant vendor's operating system documentation for all supported cancellation points. Additionally, POSIX threads defines the following function:

```
void pthread_testcancel();
```

This function may be used to test whether cancellation is required at any point within an application. It is particularly useful for allowing cancellation from a hard loop without any system cancellation points.

Asynchronous cancellation allows the thread to be canceled at any time without encountering a cancellation point. Use of asynchronous cancellation is not encouraged as it may potentially leave resources in an unknown state.

Once cancellation has been enabled within a thread, the following functions may be used to define and place a “clean up function” on the cancellation stack:

```
void pthread_cleanup_push(void (*handler, void *), void *arg);  
  
void pthread_cleanup_pop(int execute);
```

Note, in order for the application to compile, use of *pthread_cleanup_push()* requires that a matching *pthread_cleanup_pop()* is placed within the same caller function.

If multiple cleanup handlers are placed upon the cleanup stack, they will be called in LIFO order once a cancellation request is received. In order to send a cancellation request, the following function is defined:

```
int pthread_cancel(pthread_t thread);
```


The `pthread_cancel()` function returns 0 upon success or ESRCH if no such thread exists. The following code provides an example of deferred cancellation:

```

_____ cancel.c _____
1 #include <pthread.h>
2 #include <stdio.h>
3 #include <math.h>
4 #include <string.h>
5 #include <errno.h>
6 #include <unistd.h>
7 #include <stdlib.h>
8
9 double      *current_fn;
10 pthread_t   t1, t2;
11
12 void calc_fn_cleanup(void) {
13     printf("[thread %d] cancelled; cleaning up now\n",
14           (int) pthread_self());
15     if (current_fn != NULL) {
16         free(current_fn);
17     }
18     return;
19 }
20
21 void *calc_fn(void *arg) {
22     const double k = 1/sqrt(5.0), p = ((1.0+sqrt(5.0))/2.0);
23     double x;
24     int i = 1;
```

```

25     int *total = arg;
26
27     pthread_setcancelstate(PTHREAD_CANCEL_ENABLE, NULL);
28     pthread_setcanceltype(PTHREAD_CANCEL_DEFERRED, NULL);
29     pthread_cleanup_push((void *) calc_fn_cleanup, NULL);
30
31     current_fn = (double *) malloc(sizeof(double));
32
33     while (1) {
34         x = k * pow(p, ++i);
35         *current_fn = x;
36         pthread_testcancel();
37         *total = *total + 1;
38     }
39
40     pthread_cleanup_pop(1);
41     return NULL;
42 }
43
44 void *cancel_q(void) {
45     char buf[80];
46
47     while(strcasecmp(buf, "y\n") != 0) {
48         printf("[thread %d] shall I cancel [y/n]: ",
49             (int) pthread_self());
50         fgets(buf, 79, stdin);
51     }
52

```

```

53         pthread_cancel(t1);
54         return NULL;
55     }
56
57     int main(void) {
58         int     status;
59         void    *ret;
60         int     total = 0;
61
62         #ifdef __sun__
63         pthread_setconcurrency( sysconf(_SC_NPROCESSORS_ONLN)+1 );
64         #endif
65
66         status = pthread_create(&t1, NULL, (void *) calc_fn, (void *) &total);
67
68         if (status != 0) {
69             fprintf(stderr, "[thread %d] error creating thread: %s\n",
70                     pthread_self(), (char *) strerror(status));
71             exit(0);
72         }
73
74         status = pthread_create(&t2, NULL, (void *) cancel_q, NULL);
75
76         if (status != 0) {
77             fprintf(stderr, "[thread %d] error creating thread: %s\n",
78                     pthread_self(), (char *) strerror(status));
79             exit(0);
80         }

```

```

81
82     pthread_join(t1, &ret);
83
84     if (ret == PTHREAD_CANCELED) {
85         printf("[thread %d] thread %d cancelled; %d ops performed\n",
86             pthread_self(), t1, total);
87     }
88
89     return 0;
90 }

```

cancel.c

5.8 Forking and Threads

The use of *fork(2)* within a threaded application should be avoided unless it is used only to invoke another program with *exec(2)*. The POSIX standard states that only the thread which calls *fork(2)* will exist in the child. Additionally, it will own all the mutexes and have the same values for thread specific data as it did before the *fork(2)*. Additionally, to assist in performing cleanup operations before a fork, POSIX threads provide the following function (which is analogous to *atexit(2)*):

```

int pthread_atfork(void (*prepare)(void), void (*parent)(void),
                 void (*child)(void));

```

The function returns 0 if successful or ENOMEM if insufficient memory is available. The *pthread_atfork()* function allows a specified function to be

called within a thread prior to *fork(2)* being called. This allows a thread to unlock any required mutexes and protect itself from a deadlock situation in the child. All threaded software which executes non-trivial code within a child process should use *pthread_atfork()*.

Linux diverges from the POSIX standard in regards to the behavior of *fork(2)* within a threaded program. A portable method to use *fork(2)* safely on Linux and other operating systems is to create a “fork thread” which listens for requests and calls *fork(2)* and *wait(2)* on behalf of other threads.

5.9 Using Non-thread Safe Functions

Whilst thread safe versions of many standard POSIX functions (such as *strtok_r()*) have been added with the advent of threads, many non-thread safe functions still exist either in third party libraries or within libc. Functions such as *bufsplit()* can still be used within a threaded program but must be protected by a mutex. The following code provides an example of how such a function may be used:

```
...
pthread_mutex_t bufsplit_mutex = PTHREAD_MUTEX_INITIALIZER;
char **st[3];
char *buf;
```

```

buf = strdup("root:x:0:0::/sbin/sh");

pthread_mutex_lock(&bufsplit_mutex);
bufsplit(":", 0, NULL);
bufsplit(buf, (size_t) 2, st);
pthread_mutex_unlock(&bufsplit_mutex);
...

```

Third-party libraries should have their function interfaces modified accordingly to follow the standard of functions such as *strtok_r()* where any persistent data is passed as an argument to the function. However, in many cases it may not be possible to alter the interface a function presents to the outside world. In such a situation POSIX thread specific data may be of assistance. Essentially, thread specific data allows the same data structure to be used by multiple threads but each thread has it's own instance of the data and its own associated value. However, due to it's complex nature, detailed discussion of thread specific data is beyond the scope of this paper.

5.10 Advanced Threading Topics Not Covered

Much of the advanced multi-threading functionality provided within the POSIX standard has not been discussed within this paper. Specifically, thread, mutex and condition variable attributes provide the ability to control how such data types behave and what functionality they provide (e.g mutexes shared across multiple processes, etc). The standard also provides real

time extensions that provide functionality to alter scheduling characteristics of an application and its threads. Further information can be sought from threading texts, vendor documentation or the standard itself.

Additionally, The X/Open Group have provided extensions to the POSIX thread standard as part of their UNIX98 specification. The most noteworthy extension is the inclusion of read/write locks with the *pthread_rwlock_**() set of functions. These provide mutex like functionality whilst allowing a distinction between reader and writer threads thus allowing many threads to concurrently access a resource whilst allowing only a single thread to perform updates.

6 Acknowledgments

The author would like to thank the following people for reviewing and providing amendments to this paper:

David Luyer
Pacific Internet, Australia

John Ferlito
Bulletproof Networks, Australia

Nic Grant
Tripfinder Limited, Ireland

Philip Dell
Hewlett-Packard Consulting, Australia

References

- [1] Chris DiBona, Sam Ockman and Mark Stone (editors)
Open Sources - Voices from the Open Source Revolution
O'Reilly and Associates, 1999
- [2] David R. Butenhof
Programming with POSIX Threads
Addison-Wesley, 1999
- [3] E. W. Dijkstra, F. Genuys (editor)
Cooperating sequential processes, Programming languages
Academic Press, 1968
- [4] Brian O'Sullivan
comp.os.research FAQ
<http://www.serpentine.com/bos/os-faq/>, 1996
- [5] Uresh Vahalia
UNIX Internals - The New Frontiers
Prentice Hall, 1996