



BITS Pilani
Pilani Campus

Network Programming

K Hari Babu
Department of Computer Science & Information Systems



BITS Pilani
Pilani Campus



Outline

- Threads
 - Overview
 - Creation
 - Termination
 - Join
 - Detach
- Thread synchronization
 - Mutexes
 - Condition variables
- Server Design
- Multi Process Concurrency
 - Preforking models
 - Prethreading models
- Single Process Concurrency
 - Signal driven I/O
 - `epoll()`
- Concurrency in UDP

Using Mutex



- A mutex is a variable of the type `pthread_mutex_t`. Mutex must always be initialized.
 - For a statically allocated mutex

```
pthread_mutex_t mtx = PTHREAD_MUTEX_INITIALIZER;
```

```
3  #include <pthread.h>
4  int pthread_mutex_lock(pthread_mutex_t * mutex );
5  int pthread_mutex_unlock(pthread_mutex_t * mutex );
6  //Both return 0 on success, or a positive error number on error
```

- The `pthread_mutex_trylock()` function is the same as `pthread_mutex_lock()`, except that if the mutex is currently locked, `pthread_mutex_trylock()` fails, returning the error `EBUSY`.

Using Mutexes



```
1  #define NLOOP 5000
2  int      counter;          /* incremented by threads */
3  pthread_mutex_t counter_mutex = PTHREAD_MUTEX_INITIALIZER;
4  void      *doit(void *);
5  int main(int argc, char **argv)
6  {
7      pthread_t tidA, tidB;
8      Pthread_create(&tidA, NULL, &doit, NULL);
9      Pthread_create(&tidB, NULL, &doit, NULL);
10     /* wait for both threads to terminate */
11     Pthread_join(tidA, NULL);
12     Pthread_join(tidB, NULL);
13     exit(0);
14 }

15 void *
16 doit(void *vptr)
17 {
18     int      i, val;
19     for (i = 0; i < NLOOP; i++) {
20         pthread_mutex_lock(&counter_mutex);
21         val = counter;
22         printf("%d: %d\n", pthread_self(), val + 1);
23         counter = val + 1;
24         pthread_mutex_unlock(&counter_mutex);
25     }
26     return (NULL);
27 }
```

Condition Variables



- A mutex is fine to prevent simultaneous access to a shared variable, but we need something else to let us go to sleep waiting for some condition to occur.

```
1 static pthread_mutex_t mtx = PTHREAD_MUTEX_INITIALIZER;
2 static int avail = 0;
3 /*producer thread*/
4 pthread_mutex_lock(&mtx);
5 avail++;/* Let consumer know another unit is available */
6 pthread_mutex_unlock(&mtx);
7 /*consumer thread*/
8 for (;;) {
9     pthread_mutex_lock(&mtx);
10     while (avail > 0) {/* Consume all available units */
11         avail--;
12     }
13     pthread_mutex_unlock(&mtx);
14 }
```

- The above code works, but it wastes CPU time, because the consumer thread continually loops, checking the state of the variable *avail*. A condition variable remedies this problem.

Condition Variables



- Condition variable allows a thread to sleep (wait) until another thread notifies (signals) it that it must do something.
- A condition variable is always used in conjunction with a mutex.
- The mutex provides mutual exclusion for accessing the shared variable, while the condition variable is used to signal changes in the variable's state.

```
pthread_cond_t cond = PTHREAD_COND_INITIALIZER;
```

```
2  #include <pthread.h>
3  int pthread_cond_signal(pthread_cond_t * cond );
4  int pthread_cond_broadcast(pthread_cond_t * cond );
5  int pthread_cond_wait(pthread_cond_t * cond , pthread_mutex_t * mutex );
6  //All return 0 on success, or a positive error number on error
```

- Broadcast wakes up all blocked threads. Each will go through the code. Used when there is different tasks done for a particular condition.

Using Condition Variables



- Why mutex is associated with condition variable?
 - The thread locks the mutex in preparation for checking the state of the shared variable.
 - The state of the shared variable is checked.
 - If the shared variable is not in the desired state, then the thread must unlock the mutex (so that other threads can access the shared variable) before it goes to sleep on the condition variable.
 - Done atomically
 - When the thread is reawakened because the condition variable has been signaled, the mutex must once more be locked, since, typically, the thread then immediately accesses the shared variable.
- it is not possible for some other thread to acquire the mutex and signal the condition variable before the thread calling `pthread_cond_wait()` has blocked on the condition variable.

Using Condition Variables



```
1 static pthread_mutex_t mtx = PTHREAD_MUTEX_INITIALIZER;
2 static pthread_cond_t cond = PTHREAD_COND_INITIALIZER;
3 static int avail = 0;
4 /*producer thread*/
5 pthread_mutex_lock(&mtx);
6 avail++;/*Let consumer know another unit is available */
7 pthread_mutex_unlock(&mtx);
8 pthread_cond_signal(&cond); /* Wake sleeping consumer */
```

```
9 /*consumer thread*/
10 for (;;) {
11     s = pthread_mutex_lock(&mtx);
12     while (avail == 0) {/* Wait for something to consume */
13         s = pthread_cond_wait(&cond, &mtx);
14     }
15     while (avail > 0) {/* Consume all available units */
16         /* Do something with produced unit */
17         avail--;
18     }
19     s = pthread_mutex_unlock(&mtx);
20 }
```



Web client: Using threads

Threads : web client



- The client establishes an HTTP connection with a Web server and fetches a home page.
- On that page are often numerous references to other Web pages.
- Instead of fetching these other pages serially, one at a time, the client can fetch more than one at the same time, using multiple connections, one per thread.

Threads : web client



- We have designed a web client handling multiple simultaneous connections using non-blocking connect.
- Now we will design using threads:
 - With threads, we can leave the sockets in their default blocking mode.
 - Create one thread per connection.
 - Each thread can block in its call to connect. Kernel will schedule threads that are ready.

Threads : web client



- This program will read up to 20 files from a Web server.
- We specify as command-line arguments
 - the maximum number of parallel connections,
 - the server's hostname, and
 - each of the filenames to fetch from the server.

```
2 bash$ web 3 www.foobar.com image1.gif image2.gif image3.gif image4.gif  
3 image5.gif image6.gif image7.gif
```

- It means
 - three simultaneous connection
 - server's hostname
 - filename for the home page
 - the files to be read
- T1: 26.6 & 26.9

Header file



```
1  /* include web1 */
2  #include    "unpthread.h"
3  #include    <thread.h>      /* Solaris threads */
4  #define MAXFILES    20
5  #define SERV        "80"    /* port number or service name */
6  struct file {
7      char    *f_name;        /* filename */
8      char    *f_host;        /* hostname or IP address */
9      int     f_fd;           /* descriptor */
10     int     f_flags;         /* F_xxx below */
11     pthread_t f_tid;         /* thread ID */
12 } file[MAXFILES];
13 #define F_CONNECTING    1    /* connect() in progress */
14 #define F_READING       2    /* connect() complete; now reading */
15 #define F_DONE          4    /* all done */
16 #define GET_CMD         "GET %s HTTP/1.0\r\n\r\n"
17 int     nconn, nfiles, nlefttoconn, nlefttoread;
18 void    *do_get_read(void *);
19 void    home_page(const char *, const char *);
20 void    write_get_cmd(struct file *);
```

Main function

innovate

achieve

lead

```
22  int main(int argc, char **argv)
23  {
24      int          i, n, maxnconn;
25      pthread_t    tid;
26      struct file *fptr;
27      if (argc < 5)
28          err_quit("usage: web <#conns> <IPaddr> <homepage> file1 ...");
29      maxnconn = atoi(argv[1]);
30      nfiles = min(argc - 4, MAXFILES);
31      for (i = 0; i < nfiles; i++) {
32          file[i].f_name = argv[i + 4];
33          file[i].f_host = argv[2];
34          file[i].f_flags = 0;
35      }
36      printf("nfiles = %d\n", nfiles);
37      home_page(argv[2], argv[3]); /*get the homepage*/
38      nlefttoread = nlefttoconn = nfiles;
39      nconn = 0;
40      /* end web1 */
```

- Initialize structures

Main function



```
41  /* include web2 */
42  while (nlefttoread > 0) {
43      while (nconn < maxnconn && nlefttoconn > 0) {
44          /* find a file to read */
45          for (i = 0 ; i < nfiles; i++)
46              if (file[i].f_flags == 0)
47                  break;
48          file[i].f_flags = F_CONNECTING;
49          /*create a new thread*/
50          pthread_create(&tid, NULL, &do_get_read, &file[i]);
51          file[i].f_tid = tid;
52          nconn++;
53          nlefttoconn--;
54      }
55      if ( (n = pthread_join(tid, (void **) &fptr)) != 0)
56          errno = n, err sys("thr join error");
57      nconn--;
58      nlefttoread--;
59      printf("thread id %d for %s done\n", tid, fptr->f_name);
60  }
61  exit(0);
62  }
63  /* end web2 */
```

- Create maximum of maxconn threads and wait for them to terminate.

do_get_read function

innovate

achieve

lead

```
65  /* include do_get_read */
66  void *
67  do_get_read(void *vptr)
68  {
69      int          fd, n;
70      char         line[MAXLINE];
71      struct file  *fptr;
72      fptr = (struct file *) vptr;
73      fd = Tcp_connect(fptr->f_host, SERV);
74      fptr->f_fd = fd;
75      printf("do_get_read for %s, fd %d, thread %d\n",
76            fptr->f_name, fd, fptr->f_tid);
77      write_get_cmd(fptr);    /* write() the GET command */
78      /* Read server's reply */
79      for ( ; ; ) {
80          if ( (n = Read(fd, line, MAXLINE)) == 0)
81              break;          /* server closed connection */
82          printf("read %d bytes from %s\n", n, fptr->f_name);
83      }
84      printf("end-of-file on %s\n", fptr->f_name);
85      Close(fd);
86      fptr->f_flags = F_DONE;   /* clears F_READING */
87      return(fptr);           /* terminate thread */
88  }
89  /* end do_get_read */
```

Polling for Available Threads

innovate

achieve

lead

```
54 while (nlefttoread > 0) {
55 while (nconn < maxnconn && nlefttoconn > 0) {
56     for (i = 0 ; i < nfiles; i++)
57         if (file[i].f_flags == 0) break;
58     if ( (n = pthread_create(&tid, NULL, &do_get_read, &file[i])) != 0)
59         errno = n, err_sys("pthread_create error");
60     file[i].f_tid = tid; file[i].f_flags = F_CONNECTING;
61     nconn++; nlefttoconn--;
62 }
63     /* See if one of the threads is done */
64     if ( (n = pthread_mutex_lock(&ndone_mutex)) != 0)
65         errno = n, err_sys("pthread_mutex_lock error");
66     if (ndone > 0) {
67         for (i = 0; i < nfiles; i++) {
68             if (file[i].f_flags & F_DONE) {
69                 if ( (n = pthread_join(file[i].f_tid, (void **) &fptr)) != 0)
70                     errno = n, err_sys("pthread_join error");
71                 if (&file[i] != fptr)
72                     err_quit("file[i] != fptr");
73                 fptr->f_flags = F_JOINED; /* clears F_DONE */
74                 ndone--; nconn--; nlefttoread--;
75             } } }
76     if ( (n = pthread_mutex_unlock(&ndone_mutex)) != 0)
77         errno = n, err_sys("pthread_mutex_unlock error");
78 }
```



```
83 void *do_get_read(void *vptr)
84 {
85     int          fd, n;
86     char         line[MAXLINE];
87     struct file   *fptr;
88     fptr = (struct file *) vptr;
89     fd = Tcp_connect(fptr->f_host, SERV);
90     fptr->f_fd = fd;
91     printf("do_get_read for %s, fd %d, thread %d\n",
92           fptr->f_name, fd, fptr->f_tid);
93     write_get_cmd(fptr); /* write() the GET command */
94     /* Read server's reply */
95     for ( ; ; ) {
96         if ( (n = read(fd, line, MAXLINE)) <= 0) {
97             if (n == 0)
98                 break; /* server closed connection */
99             else
100                 err_sys("read error");
101         }
102         printf("read %d bytes from %s\n", n, fptr->f_name);
103     }
104     printf("end-of-file on %s\n", fptr->f_name);
105     close(fd); fptr->f_flags = F_DONE; /* clears F_READING */
106     if ( (n = pthread_mutex_lock(&ndone_mutex)) != 0)
107         errno = n, err_sys("pthread_mutex_lock error");
108     ndone++;
109     if ( (n = pthread_mutex_unlock(&ndone_mutex)) != 0)
110         errno = n, err_sys("pthread_mutex_unlock error");
111     return(fptr); /* terminate thread */
112 }
```

Using Condition Variables

innovate

achieve

lead

```
55 while (nlefttoread > 0) {
56     while (nconn < maxnconn && nlefttoconn > 0) {
57         for (i = 0 ; i < nfiles; i++)
58             if (file[i].f_flags == 0) break;
59         file[i].f_flags = F_CONNECTING;
60         Pthread_create(&tid, NULL, &do_get_read, &file[i]);
61         file[i].f_tid = tid;
62         nconn++;
63         nlefttoconn--;
64         pthread_mutex_lock(&ndone_mutex);
65         while (ndone == 0)
66             pthread_cond_wait(&ndone_cond, &ndone_mutex);
67         for (i = 0; i < nfiles; i++) {
68             if (file[i].f_flags & F_DONE) {
69                 pthread_join(file[i].f_tid, (void **) &fptr);
70                 if (&file[i] != fptr)
71                     err_quit("file[i] != fptr");
72                 fptr->f_flags = F_JOINED; /* clears F_DONE */
73                 ndone--;
74                 nconn--;
75                 nlefttoread--;
76                 printf("thread %d for %s done\n", fptr->f_tid, f
77             }
78         }
79         pthread_mutex_unlock(&ndone_mutex);
80     }
```

Using Condition Variables



```
86 void *
87 do_get_read(void *vptr)
88 {
89     fptr = (struct file *) vptr;
90     fd = Tcp_connect(fptr->f_host, SERV);
91     fptr->f_fd = fd;
92     write_get_cmd(fptr);    /* write() the GET command */
93     for ( ; ; ) {
94         if ( (n = Read(fd, line, MAXLINE)) == 0)
95             break;        /* server closed connection */
96     }
97     close(fd);
98     fptr->f_flags = F_DONE;    /* clears F_READING */
99     pthread_mutex_lock(&ndone_mutex);
100     ndone++;
101     pthread_cond_signal(&ndone_cond);
102     pthread_mutex_unlock(&ndone_mutex);
103     return(fptr);    /* terminate thread */
104 }
```

- Table shows the clock time required to fetch a Web server's home page, followed by nine image files from that server.
 - The RTT to the server is about 150 ms.
 - The home page size was 4,017 bytes and the average size of the 9 image files was 1,621 bytes.
 - TCP's segment size was 512 bytes.
- Most of the improvement is obtained with three simultaneous connections.

# simultaneous connections	Clock time (seconds), non blocking	Clock time(sec s) Threads
1	6.0	6.3
2	4.1	4.2
3	3.0	3.1
4	2.8	3.0
5	2.5	2.7
6	2.4	2.5
7	2.3	2.3
8	2.2	2.3
9	2.0	2.2



Server Design Alternatives

T1: ch30

Iterative vs Concurrent



- Iterative Servers

- Process one client at a time. Clients will experience significant delays in response.
- Suitable when:
 - suitable only when client requests can be handled quickly, since each client must wait until all of the preceding clients have been serviced.
 - A typical scenario: client and server exchange a single request and response.

- Concurrent Servers

- Handle multiple clients simultaneously. Can take advantage of CPU-I/O overlap.
- Suitable when:
 - Concurrent servers are suitable when a significant amount of processing time is required to handle each request
 - Or where the client and server engage in an extended conversation, passing messages back and forth.

Stateless Vs. Stateful Servers



- Information that a server maintains about the status of ongoing interactions with clients is called state information.
- Servers that do not keep any state information are called stateless servers; others are called stateful servers.
- Stateful Servers:
 - Keeping a small amount of information in a server
 - can reduce the size of messages
 - can allow the server to respond to requests quickly.
 - Server can compute an incremental response as each new request arrives.
- State information in a server can become incorrect if
 - messages are lost, duplicated, or delivered out of order, or if the client computer crashes and reboots.

General Server Categories



iterative connectionless	iterative connection-oriented
concurrent connectionless	concurrent connection-oriented

General Server Categories



- Iterative Connectionless
 - Common form of connectionless server.
 - Often stateless.
 - Used when it requires trivial amount of processing for each request.
- Iterative, Connection-Oriented Server
 - A less common server type
 - used for services that require a trivial amount of processing for each request, but for which reliable transport is necessary.
 - Because the overhead associated with establishing and terminating connections can be high, the average response time can be non-trivial.

General Server Categories



- Concurrent, Connectionless Server
 - An uncommon type
 - The server creates a new process to handle each request.
 - On many systems, the added cost of process creation dominates the added efficiency gained from concurrency.
 - To justify concurrency,
 - either the time required to create a new process must be significantly less than the time required to compute a response
 - or concurrent requests must be able to use many I/O devices simultaneously.

General Server Categories



- Concurrent, Connection-Oriented Server
 - The most general type of server
 - it offers reliable transport (i.e., it can be used across a wide area internet) as well as the ability to handle multiple requests concurrently.
 - Two basic implementations exist
 - concurrent processes/threads to handle multiple connections.
 - a single process and asynchronous I/O to handle multiple connections.

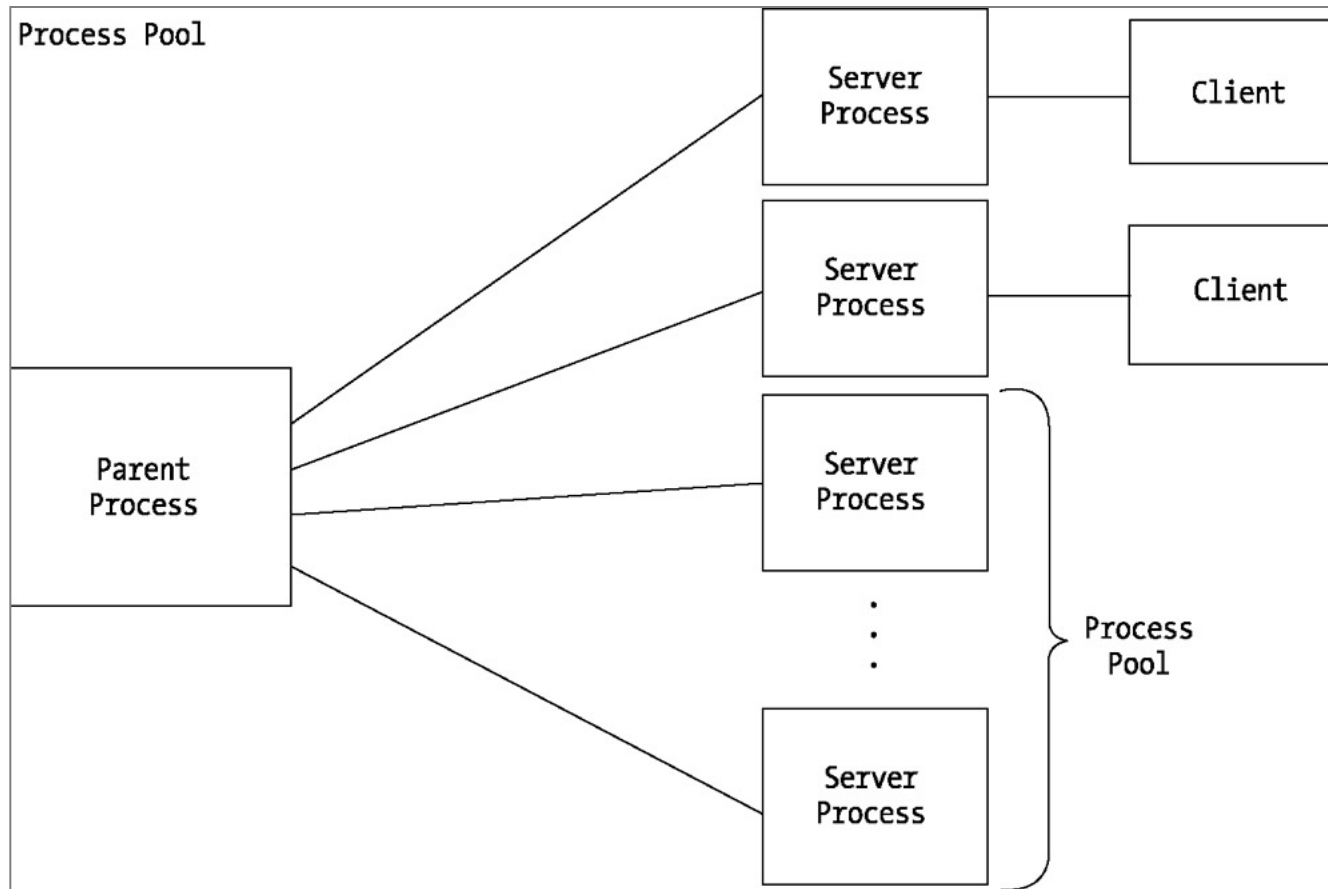
Preforked and prethreaded servers



- Traditional concurrent server model:
 - Fork a child after accepting a new client connection.
 - Good enough for low traffic services.
- For very high-load servers
 - web servers handling thousands of requests per minute
 - the cost of creating a new child (or even thread) for each client imposes a significant burden on the server.
- Instead of creating a new child process (or thread) for each client, the server precreates a fixed number of child processes (or threads) on startup.
 - Each child (thread) handles a new client. After completing one client, it accepts another connection.

- Different models
 - Child calls accept()
 - TCP Preforked Server, No Locking Around accept
 - TCP Preforked Server, Thread Locking Around accept
 - Parent calls accept() and passes the descriptor to child
 - TCP Preforked Server, Descriptor Passing

Preforking or Process Pool



Preforking or Process Pool

innovate

achieve

lead

```
static int      nchildren;
static pid_t    *pids;

int
main(int argc, char **argv)
{
    int          listenfd, i;
    socklen_t    addrlen;
    void          sig_int(int);
    pid_t        child_make(int, int, int);
    if (argc == 3)
        listenfd = Tcp_listen(NULL, argv[1], &addrlen);
    else if (argc == 4)
        listenfd = Tcp_listen(argv[1], argv[2], &addrlen);
    else
        err_quit("usage: serv02 [ <host> ] <port#> <#children>");
    nchildren = atoi(argv[argc-1]);
    pids = Calloc(nchildren, sizeof(pid_t));
    for (i = 0; i < nchildren; i++)
        pids[i] = child_make(i, listenfd, addrlen);    /* parent returns */

    Signal(SIGINT, sig_int);
    for ( ; ; )
        pause();    /* everything done by children */
}
```

Preforking or Process Pool: No Locking Around Accept

innovate

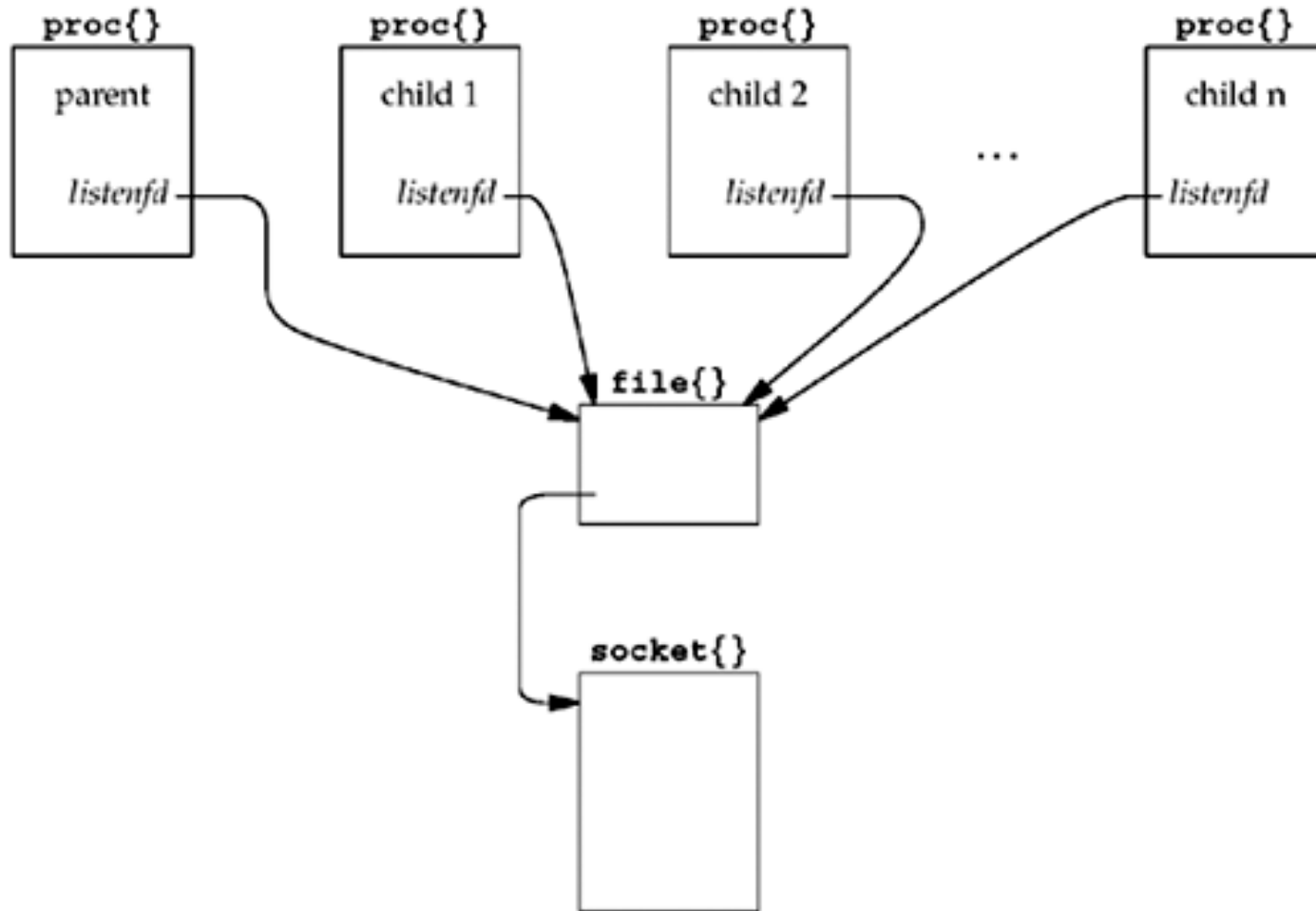
achieve

lead

```
pid_t
child_make(int i, int listenfd, int addrlen)
{
    pid_t    pid;
    void      child_main(int, int, int);
    if ( (pid = Fork()) > 0)
        return(pid);                /* parent */
    child_main(i, listenfd, addrlen); /* never returns */
}
/* end child_make */
/* include child_main */
void
child_main(int i, int listenfd, int addrlen)
{
    int                connfd;
    void               web_child(int);
    socklen_t          cliilen;
    struct sockaddr *cliaddr;
    cliaddr = Malloc(addrlen);
    printf("child %ld starting\n", (long) getpid());
    for ( ; ; ) {
        cliilen = addrlen;
        connfd = Accept(listenfd, cliaddr, &cliilen);
        web_child(connfd);           /* process the request */
        Close(connfd);
    }
}
/* end child main */
```

- Advantages:
 - No cost of fork() before responding to client.
 - Process control is simpler.
- Disadvantages:
 - Parent must guess how many children to fork.
 - If too less, clients will experience delays in response.
 - If too excessive, system performance degrades.

Thundering Herd Problem



Thundering Herd Problem



- When the program starts, N children are created, and all N call accept and all are put to sleep by the kernel.
- When the first client connection arrives, all N children are awakened.
 - because all N have gone to sleep on the same "wait channel"
because all N share the same listening descriptor.
- Even though all N are awakened, the first of the N to run will obtain the connection and the remaining $N - 1$ will all go back to sleep.

Preforking: Locking Around accept



```
static pthread_mutex_t  *mptr; /* actual mutex will be in shared memory */
void
my_lock_init(char *pathname)
{
    int          fd;
    pthread_mutexattr_t  mattr;
    fd = Open("/dev/zero", O_RDWR, 0);
    mptr = Mmap(0, sizeof(pthread_mutex_t), PROT_READ | PROT_WRITE,
                MAP_SHARED, fd, 0);

    Close(fd);
    Pthread_mutexattr_init(&mattr);
    Pthread_mutexattr_setpshared(&mattr, PTHREAD_PROCESS_SHARED);
    Pthread_mutex_init(mptr, &mattr);
}
/* end my_lock_init */
/* include my_lock_wait */
void
my_lock_wait()
{
    _Pthread_mutex_lock(mptr);
}

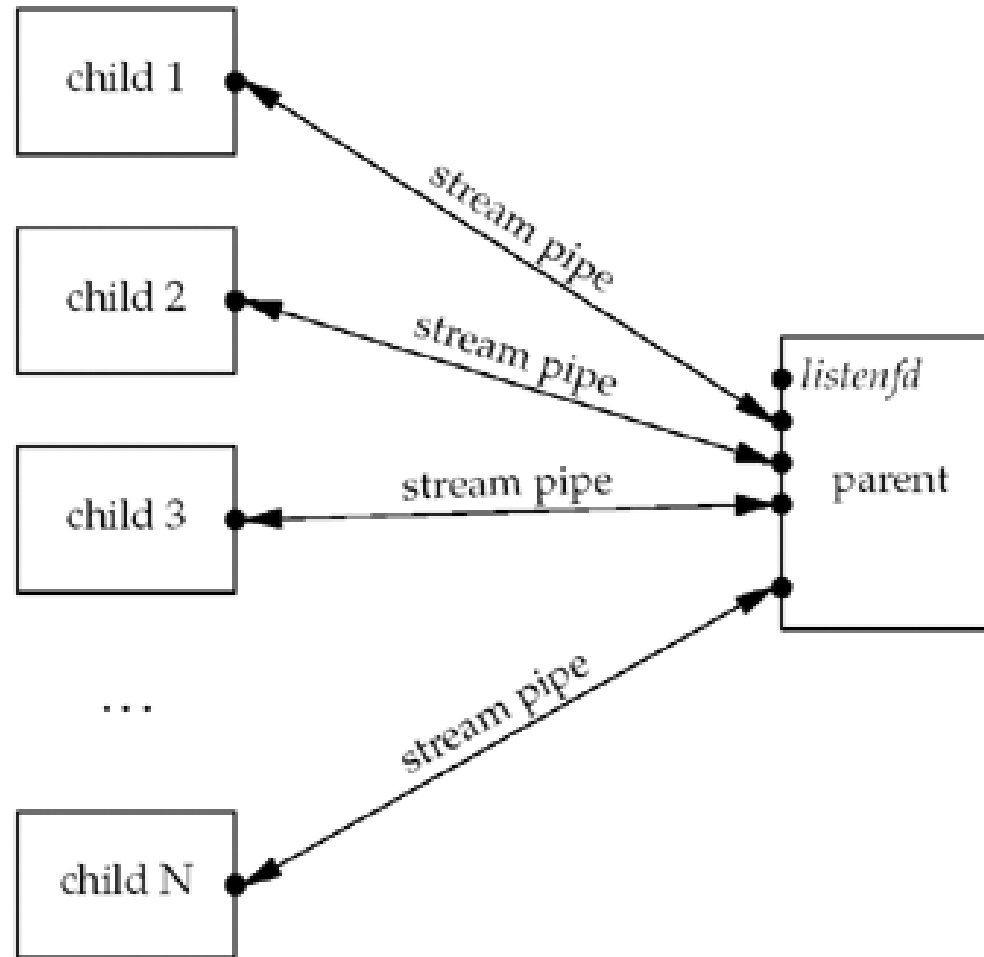
void
my_lock_release()
{
    Pthread_mutex_unlock(mptr);
}
/* end my_lock_wait */
```

Preforking: Locking Around accept



```
for ( ; ; ) {  
    clilen = addrlen;  
    +      my_lock_wait();  
    connfd = Accept(listenfd, cliaddr, &clilen);  
    +      my_lock_release();  
    web_child(connfd);          /* process request */  
    Close(connfd);
```

Preforking: Descriptor Passing



Preforking: Descriptor Passing



```
typedef struct {  
    pid_t      child_pid;           /* process ID */  
    int        child_pipefd;        /* parent's stream pipe to/from child */  
    int        child_status;        /* 0 = ready */  
    long       child_count;         /* # connections handled */  
} Child;  
  
Child *cptr;           /* array of Child structures; calloc'ed */  
~
```

- Parent maintains this structure for each child.

Preforking: Descriptor Passing

innovate

achieve

lead

```
pid_t
child_make(int i, int listenfd, int addrlen)
{
    int          sockfd[2];
    pid_t        pid;
    void          child_main(int, int, int);
    Socketpair(AF_LOCAL, SOCK_STREAM, 0, sockfd);
    if ( (pid = Fork()) > 0 ) {
        Close(sockfd[1]);
        cptr[i].child_pid = pid;
        cptr[i].child_pipefd = sockfd[0];
        cptr[i].child_status = 0;
        return(pid);          /* parent */
    }
    Dup2(sockfd[1], STDERR_FILENO);    /* child's str
    Close(sockfd[0]);
    Close(sockfd[1]);
    Close(listenfd);
    child_main(i, listenfd, addrlen);  /* never retur
}
/* end child_make */
```

Preforking: Descriptor Passing

innovate

achieve

lead

```
main(int argc, char **argv)
{
    int                listenfd, i, navail, maxfd, nsel, connfd, rc;
    void               sig_int(int);
    pid_t              child_make(int, int, int);
    ssize_t            n;
    fd_set              rset, masterset;
    socklen_t          addrlen, clilen;
    struct sockaddr     *cliaddr;

    if (argc == 3)
        listenfd = Tcp_listen(NULL, argv[1], &addrlen);
    else if (argc == 4)
        listenfd = Tcp_listen(argv[1], argv[2], &addrlen);
    else
        err_quit("usage: serv05 [ <host> ] <port#> <#children>");

    FD_ZERO(&masterset);
    FD_SET(listenfd, &masterset);
    maxfd = listenfd;
    cliaddr = Malloc(addrlen);

    nchildren = atoi(argv[argc-1]);
    navail = nchildren;
    cptr = Calloc(nchildren, sizeof(Child));

    /* 4prefork all the children */
    for (i = 0; i < nchildren; i++) {
        child_make(i, listenfd, addrlen);          /* parent returns */
        FD_SET(cptr[i].child_pipefd, &masterset);
        maxfd = max(maxfd, cptr[i].child_pipefd);
    }
}
```

Descriptor Passing

innovate

achieve

lead

```
for ( ; ; ) {
    rset = masterset;
    if (navail <= 0)
        FD_CLR(listenfd, &rset);          /* turn off if no available children */
    nsel = Select(maxfd + 1, &rset, NULL, NULL, NULL);

    /* 4check for new connections */
    if (FD_ISSET(listenfd, &rset)) {
        cliilen = addrlen;
        connfd = Accept(listenfd, cliaddr, &cliilen);

        for (i = 0; i < nchildren; i++)
            if (cptr[i].child_status == 0)
                break;                      /* available */

        if (i == nchildren)
            err_quit("no available children");
        cptr[i].child_status = 1;           /* mark child as busy */
        cptr[i].child_count++;
        navail--;

        n = Write_fd(cptr[i].child_pipefd, "", 1, connfd);
        Close(connfd);
        if (--nsel == 0)
            continue;                      /* all done with select() results */
    }
}
```

Descriptor Passing



```
    /* 4find any newly-available children */
for (i = 0; i < nchildren; i++) {
    if (FD_ISSET(cpctr[i].child_pipefd, &rset)) {
        if ( (n = Read(cpctr[i].child_pipefd, &rc, 1)) == 0)
            err_quit("child %d terminated unexpectedly", i);
        cpctr[i].child_status = 0;
        navail++;
        if (--nset == 0)
            break; /* all done with select() results */
    }
}
```



BITS Pilani
Pilani Campus



Prethreading

T1: ch 30

Preforked Server Models



- Models:
 - Parent creates pool, child calls `accept()`.
 - Parent creates pool, child calls `accept()` with a lock around.
 - Child scheduling done by kernel
 - Parent creates pool, parent calls `accept()`, parent passes connection to child.
 - Child scheduling done by parent
- Advantages
 - Robustness. Even if one child crashes, server keeps running.
 - Simple programming.
- Disadvantages
 - large context switch overheads
 - Large memory footprint per connection. Scalability issue.
 - Optimizations involving sharing information among processes (e.g., caching) harder

Prethread Server Models



- Threads have lower memory foot print and lower context switch overhead.
 - Better scalability
 - They are preferred over processes.
- Instead of creating a new thread every time, a thread pool is created on start up.
- Pthreading Server Models
 - Per-Thread accept()
 - Main thread creates thread pool, and each thread calls accept().
 - main- thread accept()
 - Main thread creates thread pool, calls accept() and pass on the connection to a thread.

Prethreaded Server per-Thread accept()



- Main thread creates *nthreads* and waits for all threads.
- Each thread calls *accept()* with mutex around.

```
1  int listenfd, nthreads;
2  socklen_t addrlen;
3  pthread_mutex_t mlock=PTHREAD_MUTEX_INITIALIZER;
4  int main(int argc, char **argv)
5  {
6      int i;
7      void sig_int(int), thread_make(int);
8      listenfd=socket();
9      bind(listenfd, );
10     nthreads = atoi(argv[argc - 1]);
11     for (i = 0; i < nthreads; i++)
12         thread_make(i); /* only main thread returns */
13     signal(SIGINT, sig_int);
14     for ( ; ; )
15         pause(); /* everything done by threads */
16 }
```

Prethreaded Server per-Thread accept()

innovate

achieve

lead

```
1 void thread_make(int i)
2 {
3     void *thread_main(void *);
4     pthread_create(&thread_tid, NULL, &thread_main, (void *) i);
5     return; /* main thread returns */
6 }
7
8 void *thread_main(void *arg)
9 {
10     int connfd;
11     void web_child(int);
12     socklen_t clilen;
13     struct sockaddr *cliaddr;
14     cliaddr = malloc(addrlen);
15     printf("thread %d starting\n", (int) arg);
16     for ( ; ; ) {
17         clilen = addrlen;
18         pthread_mutex_lock(&mlock);
19         connfd = accept(listenfd, cliaddr, &clilen);
20         pthread_mutex_unlock(&mlock);
21         tptr[(int) arg].thread_count++;
22         web_child(connfd); /* process request */
23         close(connfd);
24     }
25 }
```

Prethreaded Server Main-Thread `accept()`



- Main thread creates a pool of threads when it starts.
- Main thread calls `accept()`. Maintains a thread *tptr* array, and *clifd* array.
- *clifd* array:
 - A shared array to hold connected fds.
 - Main thread will store.
 - Child threads will take one of these.
 - *iget*: index of the next entry to be fetched by thread.
 - *iput*: index where next entry will be stored.

```
1 ▾ typedef struct {  
2     pthread_t thread_tid; /* thread ID */  
3     long      thread_count; /* # connections handled */  
4 } Thread;  
5 Thread *tptr; /* array of Thread structures; calloc'ed */  
6  
7 #define MAXNCLI 32  
8 int      clifd[MAXNCLI], iget, iput;  
9 pthread_mutex_t clifd_mutex;  
10 pthread_cond_t clifd_cond;
```

Prethreaded Server Main-Thread accept()



```
1 static int nthreads;
2 pthread_mutex_t clifd_mutex = PTHREAD_MUTEX_INITIALIZER;
3 pthread_cond_t clifd_cond = PTHREAD_COND_INITIALIZER;
4 int main(int argc, char **argv)
5 {
6     int i, listenfd, connfd;
7     void sig_int(int), thread_make(int);
8     socklen_t addrlen, clilen;
9     struct sockaddr *cliaddr;
10     listenfd=sock();
11     bind();
12     listen();
13     nthreads = atoi(argv[argc - 1]);
14     tptr = calloc(nthreads, sizeof(Thread));
15     iget = iput = 0;
16     /* create all the threads */
```

Prethreaded Server Main-Thread accept()



```
16 ▾ /* create all the threads */
17   for (i = 0; i < nthreads; i++)
18       thread_make(i);          /* only main thread returns */
19   signal(SIGINT, sig_int);
20 ▾   for ( ; ; ) {
21       clilen = addrlen;
22       connfd = accept(listenfd, cliaddr, &clilen);
23       pthread_mutex_lock(&clifd_mutex);
24       clifd[iput] = connfd;
25       if (++iput == MAXNCLI)
26           iput = 0;
27       if (iput == iget)
28           err_quit("iput = iget = %d", iput);
29       pthread_cond_signal(&clifd_cond);
30       pthread_mutex_unlock(&clifd_mutex);
31   }
32 }
```

- Condition variable *clifd_cond* is used to communicate the availability of new connection.

Prethreaded Server Main-Thread accept()



```
1 void * thread_main(void *arg)
2 {
3     int      connfd;
4     void      web_child(int);
5     printf("thread %d starting\n", (int) arg);
6     for ( ; ; ) {
7         pthread_mutex_lock(&clifd_mutex);
8         while (iget == iput)
9             pthread_cond_wait(&clifd_cond, &clifd_mutex);
10        connfd = clifd[iget]; /* connected socket to service */
11        if (++iget == MAXNCLI)
12            iget = 0;
13        pthread_mutex_unlock(&clifd_mutex);
14        tptr[(int) arg].thread_count++;
15        web_child(connfd); /* process request */
16        Close(connfd);
17    }
18 }
```

- If `iget==iput` then there is no new connection. So wait on condition variable.

Comparing Multi Process/Multi Thread Designs



- Maximum 10 simultaneous connections.
- Process pool size:15 / Thread pool size: 15

Model	Process control CPU time (secs)
Iterative	0 (base case)
One fork per client req	20.90
Prefork with child calling accept	1.80
Prefork with child calling accept mutex around	1.75
Prefork with parent passing socket fd to child	2.58
Thread per client req	0.99
Pre threaded with child calling accept	1.93
Prethreaded with main thread calling accept	2.05

- Pre threaded models are better than preforked models.
- Kernel managed connection distribution better performance.

Threads vs Processes



- Threads provide better performance than processes.
 - Lower context switch overheads
 - Shared address space simplifies optimizations (e.g., caches)
- But
 - Some extra memory needed to support multiple stacks
 - Need thread-safe programs, synchronization
 - Security: one faulty thread can bring down whole server.
- Apache combines best of both processes and threads using Preforked and prethreaded model.
- IIS on windows platform supports only multi threaded model.

- Apache is a open source HTTP web server and is built and maintained over at Apache.org
- Apache is comprised of two main building blocks
 - Apache core
 - Apache modules
- Easy to implement and easy to extend its abilities by adding different modules.
- More info at
http://www.shoshin.uwaterloo.ca/~oadragoi/cw/CS746G/a1/apache_conceptual_arch.html

Multi-Processing Modules



- Apache 1.3 is a pre-forking server.
 - Easier for UNIX platforms but difficult in Windows platform.
- In Apache 2.0, an abstract layer for Multi-Processing Modules is designed.

Concurrency in Apache

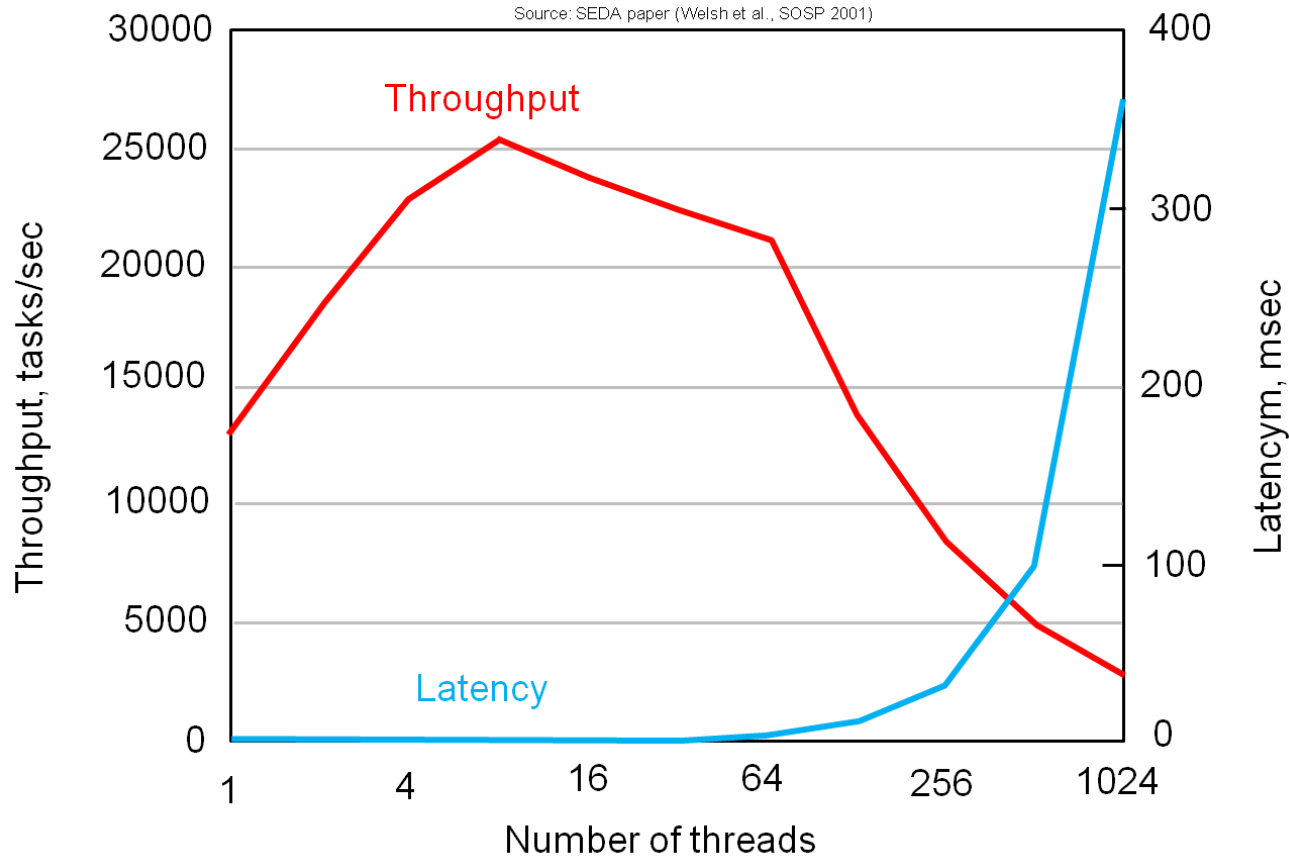


- Following models are supported by Apache
 - prefork
 - Default. Currently the default for Unix and sites that require stability.
 - threaded
 - Suitable for sites that require the benefits brought by threading, particularly reduced memory footprint and improved interthread communications.
 - mpmt_pthread
 - Similar to prefork, but each child process has a specified number of threads. It is possible to specify a minimum and maximum number of idle threads.
 - Dexter
 - Multiprocess, multithreaded MPM that allows you to specify a static number of processes.



Single Process Server Models

Threads Scalability



- As the number of threads increase in the system, more time is taken by context switching than actually doing productive work.

C10k Problem



- The C10k problem is the problem of providing scalable concurrency to handle a large number of clients at the same time.
- <http://www.kegel.com/c10k.html>
- Proposes alternatives to thread/process based servers.

Single Process Servers



- We can design a single server process to handle multiple clients employing I/O multiplexing, signal-driven I/O or epoll.
- The server process must take on some of the scheduling tasks that are normally handled by the kernel.
 - Signal-driven I/O
 - a process requests that the kernel send it a signal when input is available or data can be written on a specified file descriptor.
 - When monitoring large numbers of file descriptors, signal-driven I/O performs better than `select()` and `poll()`.
 - POSIX AIO
 - Linux provides a threads-based implementation of POSIX AIO within glibc. Not widely used.
 - epoll
 - Scalable for large number of fds. Specific to Linux. Please see R1:63.4.

- level-triggered interrupts occur whenever the file descriptor is ready for I/O
 - 1000 bytes of data in receive buffer
 - you call `recv()` and extract 500 bytes
 - `select()` will continue to indicate the fd is ready because there are still 500 bytes in the buffer
- edge-triggered interrupts occur whenever the file descriptor goes from being not ready to ready
 - 1000 bytes of data in receive buffer . Kernel delivers a signal to owner process.
 - you call `recv()` and extract 500 bytes
 - Another signal will not be delivered until the receive buffer goes down to zero and then back up to some positive number

Level-Triggered and Edge-Triggered Notification



I/O model	Level-triggered?	Edge-triggered?
<i>select()</i> , <i>poll()</i>	•	
Signal-driven I/O		•
<i>epoll</i>	•	•

- Why *epoll()* performs better?
 - On each call to *select()* or *poll()*, the kernel must check all of the file descriptors specified in the call.
 - But in *epoll* using *epoll_ctl()* fd list is created in kernel space.
 - Whenever I/O becomes ready on a fd, kernel adds it to *ready list*. When user calls *epoll_wait()*, simply return *ready list*.
 - *select()* passes data to kernel each time it is called.

Number of descriptors monitored (<i>N</i>)	<i>poll()</i> CPU time (seconds)	<i>select()</i> CPU time (seconds)	<i>epoll</i> CPU time (seconds)
10	0.61	0.73	0.41
100	2.9	3.0	0.42
1000	35	35	0.53
10000	990	930	0.66

poll()



```
2 #include <poll.h>
3 int poll(struct pollfd fds [], nfds_t nfds , int timeout );
4 /*Returns number of ready file descriptors, 0 on timeout, or -1 on error*/
```

- With select(), we provide three sets, each marked to indicate the file descriptors of interest.
- With poll(), we provide a list of file descriptors, each marked with the set of events of interest.

```
struct pollfd {
    int    fd;        /* File descriptor */
    short  events;     /* Requested events bit mask */
    short  revents;    /* Returned events bit mask */
};
```

- The caller initializes *events* to specify the events to be monitored for the file descriptor *fd*. When *poll()* returns, *revents* is set to indicate which of those events occurred for this file descriptor.
 - *events* can be 0 if do not want to include that *fd*.

Table 63-2: Bit-mask values for *events* and *revents* fields of the *pollfd* structure

Bit	Input in <i>events</i> ?	Returned in <i>revents</i> ?	Description
POLLIN	•	•	Data other than high-priority data can be read
POLLRDNORM	•	•	Equivalent to POLLIN
POLLRDBAND	•	•	Priority data can be read (unused on Linux)
POLLPRI	•	•	High-priority data can be read
POLLRDHUP	•	•	Shutdown on peer socket
POLLOUT	•	•	Normal data can be written
POLLWRNORM	•	•	Equivalent to POLLOUT
POLLWRBAND	•	•	Priority data can be written
POLLERR		•	An error has occurred
POLLHUP		•	A hangup has occurred
POLLNVAL		•	File descriptor is not open
POLLMSG			Unused on Linux (and unspecified in SUSv3)

- flags of real interest are POLLIN, POLLOUT, POLLPRI, POLLRDHUP, POLLHUP, and POLLERR.

poll() example

innovate

achieve

lead

```
2  /* Build the file descriptor list to be supplied to poll(). This list
3     is set to contain the file descriptors for the read ends of all of
4     the pipes. */
5  for (j = 0; j < numPipes; j++) {
6     pollFd[j].fd = pfd[j][0];
7     pollFd[j].events = POLLIN;
8  }
9  ready = poll(pollFd, numPipes, -1);          /* Nonblocking */
10 if (ready == -1)
11     errExit("poll");
12 printf("poll() returned: %d\n", ready);
13 /* Check which pipes have data available for reading */
14 for (j = 0; j < numPipes; j++)
15     if (pollFd[j].revents & POLLIN)
16         printf("Readable: %d %3d\n", j, pollFd[j].fd);
17 exit(EXIT_SUCCESS);
18 }
```

Differences between select() and poll()



- When fds are sparsely present, poll() gives better performance than select().
- select() is widely supported.
- Mapping between select() and poll():

```
/*To implement select(), a set of macros is used to convert the information
returned by the kernel poll routines into the corresponding event types returned
by select():*/
#define POLLIN_SET (POLLRDNORM | POLLRDBAND | POLLIN | POLLHUP | POLLERR)
/* Ready for reading */
#define POLLOUT_SET (POLLWRBAND | POLLWRNORM | POLLOUT | POLLERR)
/* Ready for writing */
#define POLLEX_SET (POLLPRI)
```

select() & poll()



- CPU time required by select() and poll() increases with the number of file descriptors being monitored.
- The poor scaling performance of select() and poll() stems from :
 - a program makes repeated calls to monitor the same set of file descriptors; however, the kernel doesn't remember the list of file descriptors to be monitored between successive calls.
- Signal-driven I/O and epoll
 - allow the kernel to record a persistent list of file descriptors
 - scale according to the number of I/O events that occur, rather than according to the number of file descriptors being monitored

Signal-Driven I/O Model



- To use signal-driven I/O with a socket (SIGIO) requires the process to perform the following three steps:
 - A signal handler must be established for the SIGIO signal.
 - The socket owner must be set, normally with the F_SETOWN command of fcntl.
 - Signal-driven I/O must be enabled for the socket, normally with the F_SETFL command of fcntl to turn on the O_ASYNC flag.

Signal-Driven I/O Model



```
1  /* Establish handler for "I/O possible" signal */
2  sigemptyset(&sa.sa_mask);
3  sa.sa_flags = SA_RESTART;
4  sa.sa_handler = sigioHandler;
5  if (sigaction(SIGIO, &sa, NULL) == -1)
6      errExit("sigaction");
7  /* Set owner process that is to receive "I/O possible" signal */
8  if (fcntl(STDIN_FILENO, F_SETOWN, getpid()) == -1)
9      errExit("fcntl(F_SETOWN)");
10 /* Enable "I/O possible" signaling and make I/O nonblocking
11    for file descriptor */
12 flags = fcntl(STDIN_FILENO, F_GETFL);
13 if (fcntl(STDIN_FILENO, F_SETFL, flags | O_ASYNC | O_NONBLOCK) == -1)
14     errExit("fcntl(F_SETFL)");
15
16 if (gotSigio) { /* Is input available? */
17     /* Read all available input until error (probably EAGAIN)
18        hash (#) character is read */
19     while (read(STDIN_FILENO, &ch, 1) > 0 && !done) {
20         printf("cnt=%d; read %c\n", cnt, ch);
21         done = ch == '#';
22     }
23     gotSigio = 0;
24 }
25 }
```

```
28 static void
29 sigioHandler(int sig)
30 {
31     gotSigio = 1;
32 }
```


Signal driven I/O for large no of FDs



- Two more steps:
- Employ `fcntl()` operation, `F_SETSIG` , to specify a realtime signal that should be delivered instead of `SIGIO` when I/O is possible on a file descriptor.
 - Realtime signals allow queuing of signals
- Specify the `SA_SIGINFO` flag when using `sigaction()` to establish the handler for the realtime signal employed in the previous step.
 - `si_signo`: the number of the signal that caused the invocation of the handler.
 - This value is the same as the first argument to the signal handler.
 - `si_fd`: the file descriptor for which the I/O event occurred.
 - `si_code`: a code indicating the type of event that occurred.
 - `si_band`: a bit mask containing the same bits as are returned in the `revents` field by the `poll()` system call.

siginfo structure



```
1  typedef struct {
2      int      si_signo;      /* Signal number */
3      int      si_code;      /* Signal code */
4      int      si_trapno;    /* Trap number for hardware-generated signal
5                               (unused on most architectures) */
6      union sigval si_value; /* Accompanying data from sigqueue() */
7      pid_t    si_pid;      /* Process ID of sending process */
8      uid_t    si_uid;      /* Real user ID of sender */
9      int      si_errno;     /* Error number (generally unused) */
10     void     *si_addr;     /* Address that generated signal
11                               (hardware-generated signals only) */
12     int      si_overrun;    /* Overrun count (Linux 2.6, POSIX timers) */
13     int      si_timerid;    /* (Kernel-internal) Timer ID
14                               (Linux 2.6, POSIX timers) */
15     long     si_band;      /* Band event (SIGPOLL/SIGIO) */
16     int      si_fd;        /* File descriptor (SIGPOLL/SIGIO) */
17     int      si_status;     /* Exit status or signal (SIGCHLD) */
18     clock_t  si_utime;     /* User CPU time (SIGCHLD) */
19     clock_t  si_stime;     /* System CPU time (SIGCHLD) */
20 } siginfo_t;
```

Table 63-7: *si_code* and *si_band* values in the *siginfo_t* structure for “I/O possible” events

<i>si_code</i>	<i>si_band</i> mask value	Description
POLL_IN	POLLIN POLLRDNORM	Input available; end-of-file condition
POLL_OUT	POLLOUT POLLWRNORM POLLWRBAND	Output possible
POLL_MSG	POLLIN POLLRDNORM POLLMSG	Input message available (unused)
POLL_ERR	POLLERR	I/O error
POLL_PRI	POLLPRI POLLRDNORM	High-priority input available
POLL_HUP	POLLHUP POLLERR	Hangup occurred

Server with Signal Driven I/O



```
90 struct sigaction sa, sa1;
91 memset (&sa, '\0', sizeof (sa));
92 memset (&sa, '\0', sizeof (sa1));
93 sigemptyset (&sa.sa_mask);
94 sa.sa_flags = SA_SIGINFO;
95 sa.sa_sigaction = &sigioListenHandler; //for accepting new conn
96 sigaction (SIGIO, &sa, NULL);
97 sigaction (SIGRTMIN + 1, &sa, NULL);
98
99 sigemptyset (&sa1.sa_mask);
100 sa1.sa_flags = SA_SIGINFO;
101 sa1.sa_sigaction = &sigioConnHandler; //for reading data
102 sigaction (SIGRTMIN + 2, &sa1, NULL);
```


- Two Realtime signals are used. One on listenfd and other on connfds.
 - Unlike standard signals, real-time signals have no predefined meanings.
 - They are queued. They should be defined as SIGRTMIN+n because SIGRTMIN may vary across OSs.

Server with Signal Driven I/O



```
104 listenfd = socket (AF_INET, SOCK_STREAM, 0);
105 bzero (&servaddr, sizeof (servaddr));
106 servaddr.sin_family = AF_INET;
107 servaddr.sin_addr.s_addr = htonl (INADDR_ANY);
108 servaddr.sin_port = htons (atoi (argv[1]));
109 bind (listenfd, (struct sockaddr *) &servaddr, sizeof (servaddr));
110 listen (listenfd, LISTENQ);
111
112 fcntl (listenfd, F_SETOWN, getpid ());
113 int flags = fcntl (listenfd, F_GETFL); /* Get current flags */
114 fcntl (listenfd, F_SETFL, flags | O_ASYNC | O_NONBLOCK); //set signal driven IO
115 fcntl (listenfd, F_SETSIG, SIGRTMIN + 1); //replace SIGIO with realtime signal
```

- Line 114: Set listening socket to receive a signal on IO availability.
- Line 115: Replace default signal SIGIO with realtime signal SIGRTMIN+1.



```

20 int listenfd; //global var so that signal handlers can access them.
21 int connfd;
22 static void
23 sigioListenHandler (int sig, siginfo_t * si, void *ucontext)
24 {
25     printf ("no:%d, for fd:%d,  event band:%ld\n", si->si_signo,
26         (int) si->si_fd, (long) si->si_band);
27     fflush (stdout);
28     if (si->si_code==POLL_IN)
29     {
30         int n = accept (listenfd, NULL, 0);
31         if (n > 0)
32             connfd = n;
33         fcntl (connfd, F_SETOWN, getpid ());
34         int flags = fcntl (connfd, F_GETFL); /* Get current flags */
35         fcntl (connfd, F_SETFL, flags | O_ASYNC | O_NONBLOCK);
36         fcntl (connfd, F_SETSIG, SIGRTMIN + 2);
37     }
38     if (sig == SIGIO)
39         printf ("Real time signalQ overflow");
40
41 }

```

- This handler is for listenfd. It receives siginfo_t when a signal is delivered.
- si_code carries the event that has occurred.
 - Accept connection and set new socket to receive SIGRTMIN+2 signal.

```
44 static void
45 sigioConnHandler (int sig, siginfo_t * si, void *ucontext)
46 {
47     printf ("no:%d, for fd:%d, , event code:%d, event band:%ld\n",
48         si->si_signo, (int) si->si_fd, (int) si->si_code,
49         (long) si->si_band);
50     fflush (stdout);
51     if (si->si_code == POLL_IN)
52     {
53         //input available
54         int n = read (si->si_fd, buf, MAXLINE);
55         if (n == 0)
56         {
57             close (si->si_fd);
58             printf ("Socket %d closed\n", si->si_fd);
59         }
60         else if (n > 0)
61         {
62             buf[n] = '\0';
63             printf ("Data from connfd %d: %s %d\n", connfd, buf, n);
64             write (si->si_fd, "OK", 2);
65         }
66     }
```

- At line 51, if the event is POLL_IN, read data from the socket and write back the data to the socket.
- If EOF is received, close the socket.

- The central data structure of the epoll API is an epoll instance.
- It serves two purposes:
 - recording a list of file descriptors that this process has declared an interest in monitoring—the interest list; and
 - maintaining a list of file descriptors that are ready for I/O—the ready list.
- The epoll API consists of three system calls:
 - The `epoll_create()` system call creates an epoll instance and returns a file descriptor.
 - The `epoll_ctl()` system call manipulates the interest list associated with an epoll. Add/del/modify a fd.
 - The `epoll_wait()` system call returns items from the ready list associated with an epoll instance.

epoll



```
1  #include <sys/epoll.h>
2  int epoll_create(int size );
3  /*Returns file descriptor on success, or -1 on error*/
```

```
5  #include <sys/epoll.h>
6  int epoll_ctl(int epfd , int op , int fd , struct epoll_event * ev );
7  //Returns 0 on success, or -1 on error
```

```
struct epoll_event {
    uint32_t      events;      /* epoll events (bit mask) */
    epoll_data_t data;        /* User data */
};

typedef union epoll_data {
    void          *ptr;        /* Pointer to user-defined data */
    int            fd;         /* File descriptor */
    uint32_t       u32;        /* 32-bit integer */
    uint64_t       u64;        /* 64-bit integer */
} epoll_data_t;
```

epoll_ctl()



- EPOLL_CTL_ADD
 - Add the file descriptor fd to the interest list for epfd.
- EPOLL_CTL_MOD
 - Modify the events setting for the file descriptor fd, using the information
- EPOLL_CTL_DEL
 - Remove the file descriptor fd from the interest list for epfd.

```
1  int epfd;
2  struct epoll_event ev;
3  epfd = epoll_create(5);
4  if (epfd == -1)
5      errExit("epoll_create");
6  ev.data.fd = fd;
7  ev.events = EPOLLIN;
8  if (epoll_ctl(epfd, EPOLL_CTL_ADD, fd, ev) == -1)
9      errExit("epoll_ctl");
```

epoll_wait()



```
#include <sys/epoll.h>
int epoll_wait(int  epfd , struct epoll_event * evlist , int  maxevents , int  timeout );
//Returns number of ready file descriptors, 0 on timeout, or -1 on error
```

milli secs
-1=blocking, 0=non-blockin
>0 timeout

- The *epoll_wait()* system call returns list of ready file descriptors of epoll instance *epfd*.
- Ready file descriptors is returned in the array of *epoll_event* structures pointed to by *evlist*.
 - Allocated by caller, *maxevents* is the no of structures in *evlist*.
 - Each structure *evlist* has information about a single ready fd.
 - The *events* subfield returns a mask of the events that have occurred on this fd.
 - The *data* subfield returns whatever value was specified in *ev.data* when we registered interest in this fd using *epoll_ctl()*.
 - data field is the only mechanism for finding out the fd.

epoll Events



Table 63-8: Bit-mask values for the *epoll events* field

Bit	Input to <i>epoll_ctl()</i> ?	Returned by <i>epoll_wait()</i> ?	Description
EPOLLIN	•	•	Data other than high-priority data can be read
EPOLLPRI	•	•	High-priority data can be read
EPOLLRDHUP	•	•	Shutdown on peer socket (since Linux 2.6.17)
EPOLLOUT	•	•	Normal data can be written
EPOLLET	•		Employ edge-triggered event notification
EPOLLONESHOT	•		Disable monitoring after event notification
EPOLLERR		•	An error has occurred
EPOLLHUP		•	A hangup has occurred

Server with epoll



```
55 epfd = epoll_create (20);
56 if (epfd == -1)
57     errExit ("epoll_create");
58 ev.events = EPOLLIN;          /* Only interested in input events */
59 ev.data.fd = listenfd;
60 if (epoll_ctl (epfd, EPOLL_CTL_ADD, listenfd, &ev) == -1)
61     errExit ("epoll_ctl");
62 for (;;)
63 {
64     ready = epoll_wait (epfd, evlist, MAX_EVENTS, -1);
65     if (ready == -1)
66     {
67         if (errno == EINTR)
68             continue;        /* Restart if interrupted by signal */
69         else
70             errExit ("epoll_wait");
71     }
```

- At line no 55, epoll instance is created. At 60, listenfd is added to interest list on event EPOLLIN.
- `epoll_wait` will block until a fd becomes available.
 - Diff between `select()` and `epoll_wait()` is: `epoll_wait` returns only available fds.

Server with epoll



```
72     for (j = 0; j < ready; j++)
73     {
74         if (evlist[j].events & EPOLLIN)
75         {
76             if (evlist[j].data.fd == listenfd)
77             {
78                 clilen = sizeof (cliaddr);
79                 char ip[128];
80                 memset (ip, '\0', 128);
81                 int connfd =
82                     accept (listenfd, (struct sockaddr *) &cliaddr, &clilen);
83                 ev.events = EPOLLIN; /* Only interested in input events */
84                 ev.data.fd = connfd;
85                 if (epoll_ctl (epfd, EPOLL_CTL_ADD, connfd, &ev) == -1)
86                     errExit ("epoll_ctl");
87             }
88         }
89     }
```

- Test all returned in evlist array.
- If listenfd is set, accept a new connection. Add new connfd to interest list.
 - Note that we do not need separate client array here. Unless we need them in for statistical purposes.

Server with epoll



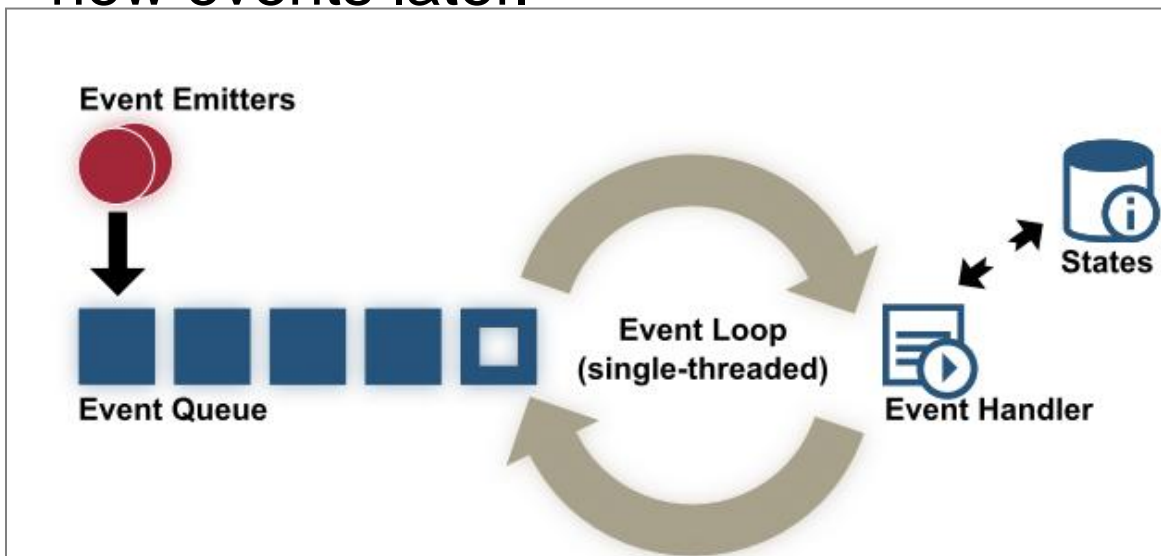
```
88     else
89     {
90         int s = read (evlist[j].data.fd, buf, MAX_BUF);
91         buf[s] = '\0';
92         if (s == -1)
93             errExit ("read");
94         if (s == 0)
95         {
96             close (evlist[j].data.fd);
97         }
98         if (s > 0)
99             write (evlist[j].data.fd, buf, strlen (buf));
100
101     }
102 }
```

- If fd is not listenfd, read data, process and send back to the client.
- If EOF is encountered, close the socket.
 - Closing a socket automatically removes it from interest list.

Event Driven Architectures



- A single threaded event loop consumes event after event from the queue and sequentially executes associated event handler code.
- New events are emitted by external sources such as socket or file I/O notifications.
- Event handlers trigger I/O actions that eventually result in new events later.



```
1 while (1) {  
2   events = getEvents();  
3   for (e in events)  
4     processEvent(e);  
5 }
```


Event Driven Architectures

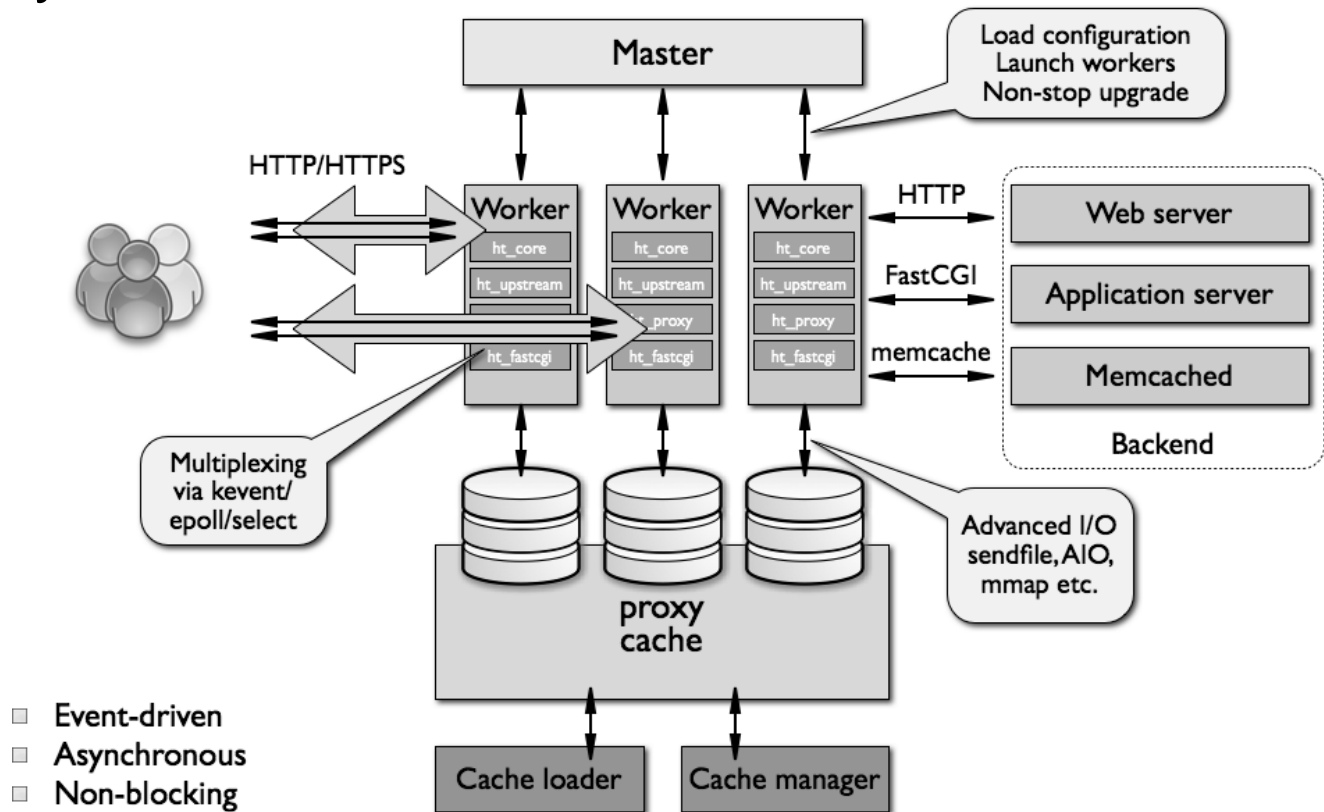


- Thread is very similar to a scheduler, multiplexing multiple connections to a single flow of execution.
- The states of the connections are organized in appropriate data structures— using finite state machines etc.
- Event-driven server architectures is dependent on the availability of asynchronous/non-blocking I/O operations at OS level.

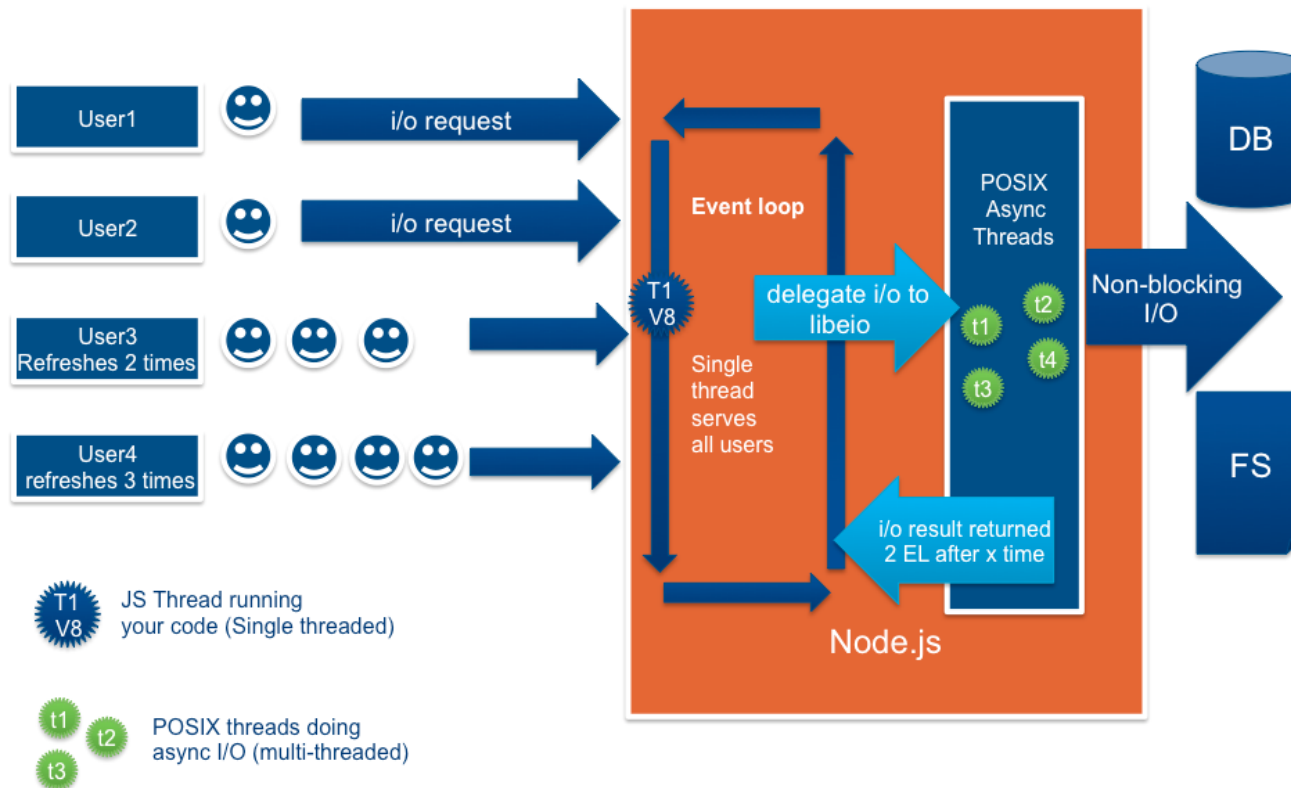
	Blocking	Non-blocking
Synchronous	read/write	read/write using O_NONBLOCK
Asynchronous	I/O multiplexing (select/poll/epoll)	AIO

- IO Multiplexing with Threads
 - Availability of data
 - Copying through a helper thread or process.
 - Notification through a call back function.
- AIO
 - Completion of data.
 - Call back functions which modify the state of the connection and generate events.

- Multiple worker processes
- Shared listening socket
- Uses Asynchronous callback functions.

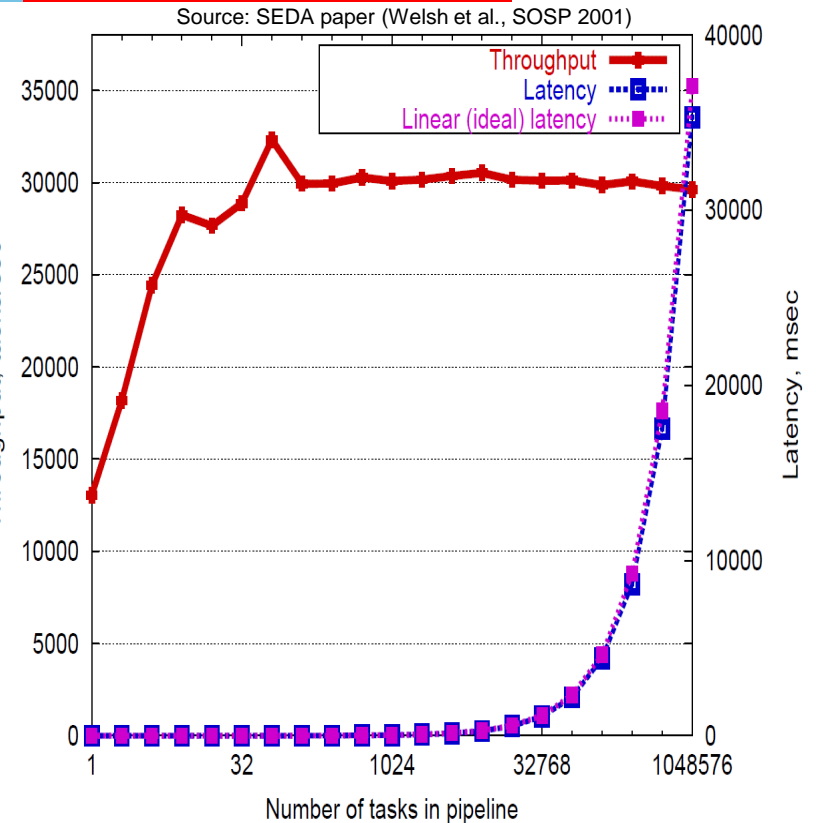
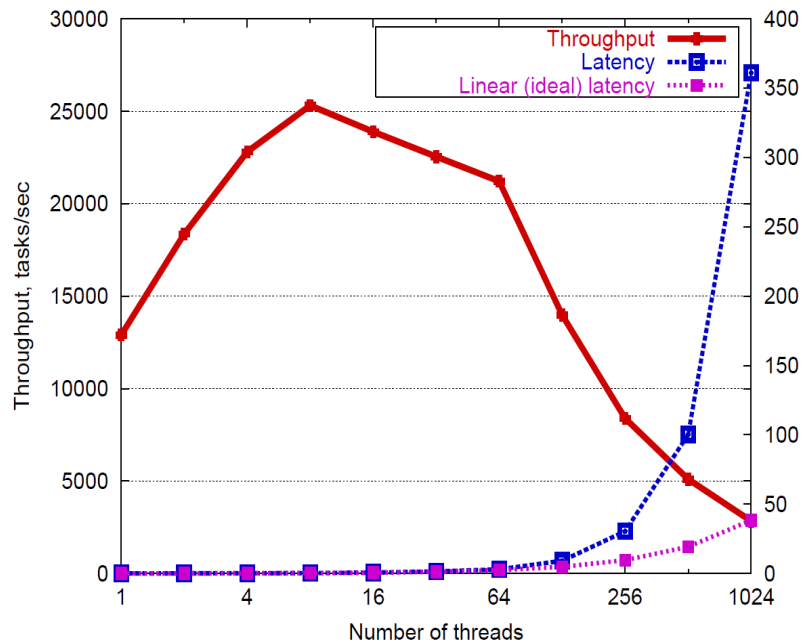


- Single thread
- Asynchronous callback functions



- Two patterns that involve event demultiplexors are called Reactor and Proactor.
 - The Reactor patterns involve synchronous I/O, whereas the Proactor pattern involves asynchronous I/O.
 - In Reactor, the event demultiplexor waits for events that indicate when a file descriptor or socket is ready for a read or write operation.
 - The demultiplexor passes this event to the appropriate handler, which is responsible for performing the actual read or write.
 - Proactor pattern, the event demultiplexor initiates asynchronous read and write operations.
 - The event demultiplexor waits for events that indicate the completion of the I/O operation, and forwards those events to the appropriate handlers.

Threads vs events



- No throughput degradation under load
- Peak throughput is higher



Concurrency UDP Servers

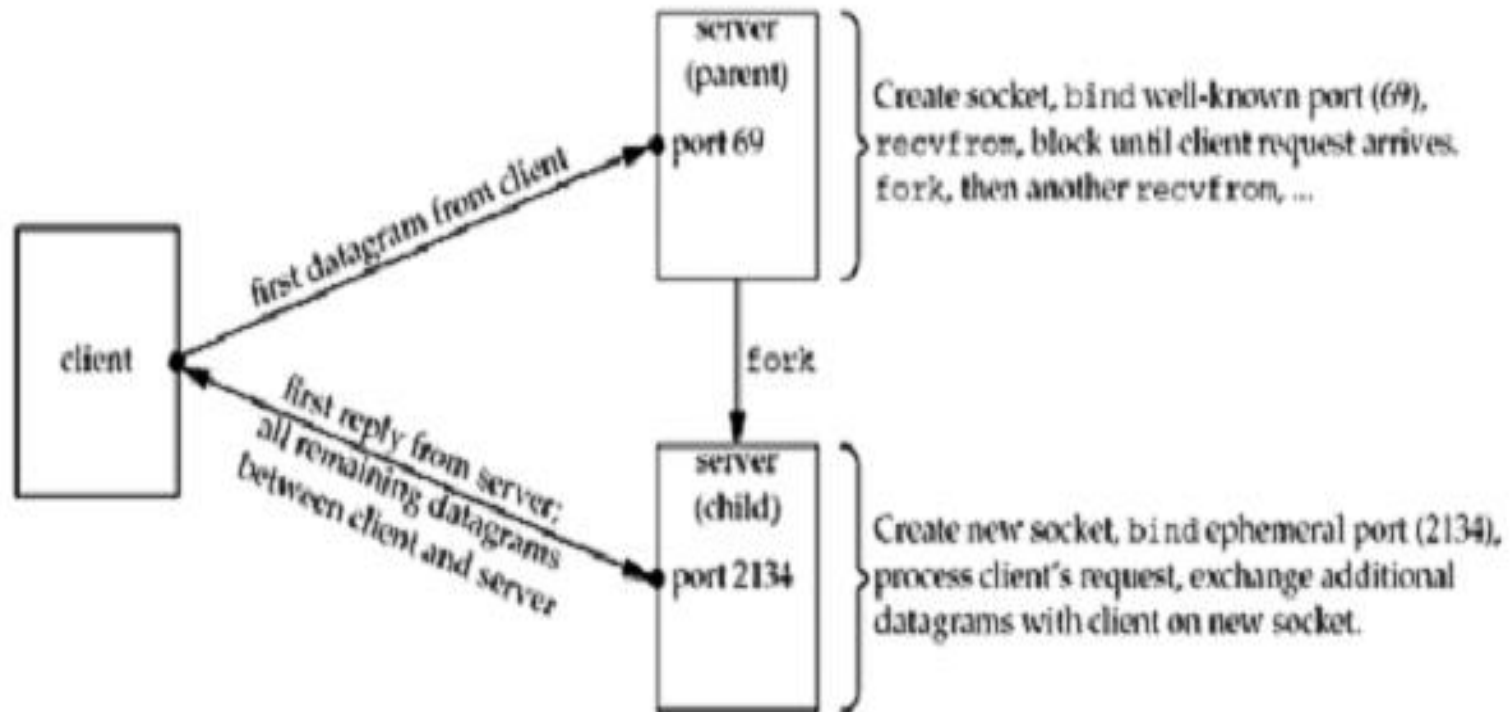
T1: ch 22.7

Concurrent UDP Servers



- Two different types of servers:
- First is a simple UDP server that reads a client request, sends a reply, and is then finished with the client
 - Concurrency: fork a child and let it handle the request
- Second is a UDP server that exchanges multiple datagrams with the client. Extended conversation.
 - Create a new socket for each client, bind an ephemeral port to that socket, and use that socket for all its replies.
 - The client looks at the port number of the server's first reply and send subsequent datagrams to that port.

Concurrency in UDP for Extended Conversations



Acknowledgements



Q&A





BITS Pilani
Pilani Campus



Thank You