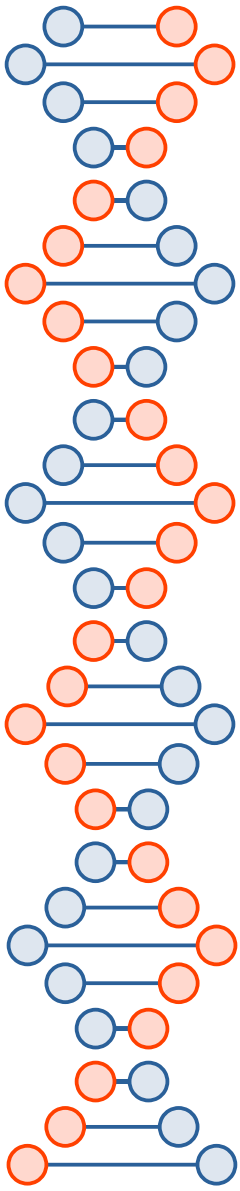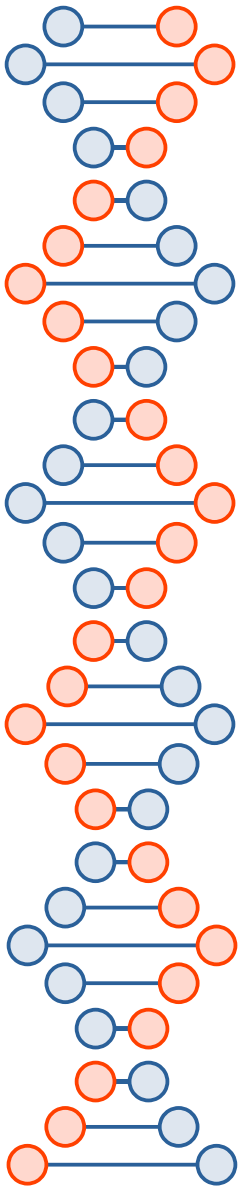# Signals and Event-Based I/O

## CSE 536 Spring 2024
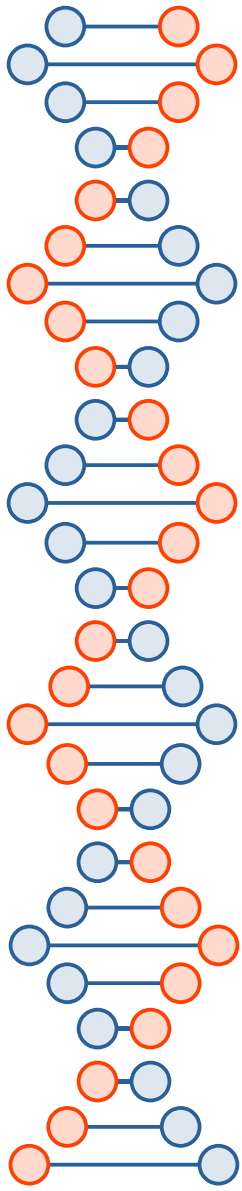jedimaestro@asu.edu

# Why Event-Based I/O?

- Multithreading can lead to a lot of errors, complexity

- Blocking is bad for performance

  - Blocking means your process is put in a wait queue because of a system call you made, basically

# Outline

- UNIX signals

- poll and ppoll()
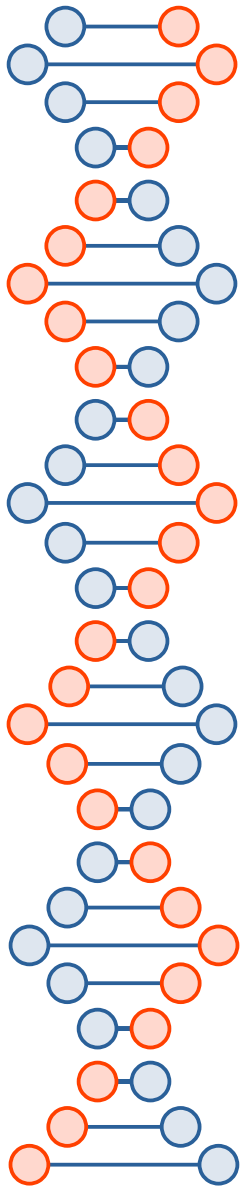
- select() and pselect()

- epoll()

- kqueue()

File    Machine    View    Input    Devices    Help

```
jedi@server:~$ cat mysignals.c
#include <stdio.h>
#include <stdlib.h>
#include <signal.h>

void sig_handler(int signum)
{
        printf("Received signal %d\n", signum);
}

void other_sig_handler(int signum)
{
        printf("Got signal %d\n", signum);
}

int main()
{
        signal(SIGINT, sig_handler);
        signal(SIGUSR1, sig_handler);
        signal(SIGUSR2, other_sig_handler);
        signal(SIGFPE, other_sig_handler);
        while(1);
        return 0;
}
jedi@server:~$ gcc mysignals.c -o mysignals
jedi@server:~$ ./mysignals
Received signal 2
Received signal 10
Got signal 12
Got signal 8
jedi@server:~$ _
```
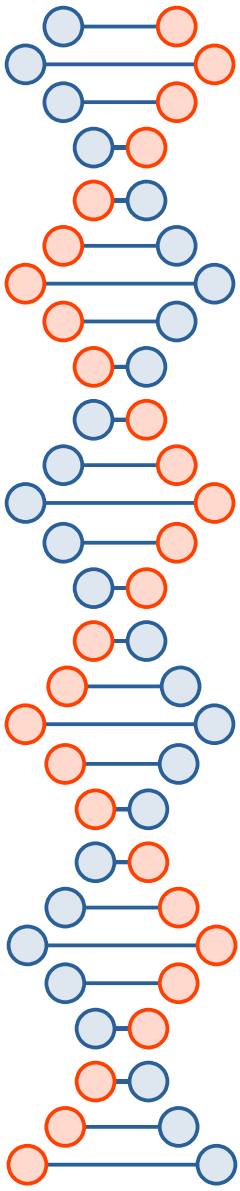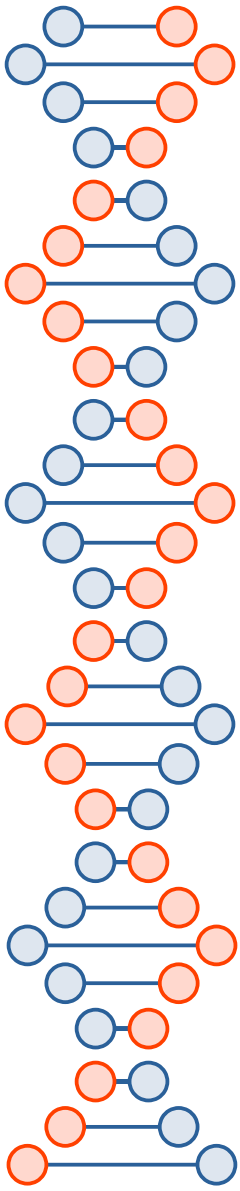
Right Ctrl

4

```
jedi@server:~$ pidof mysignals
2980
jedi@server:~$ kill -SIGINT 2980
jedi@server:~$ kill -SIGUSR1 2980
jedi@server:~$ kill -SIGUSR2 2980
jedi@server:~$ kill -SIGFPE 2980
jedi@server:~$ kill -SIGPIPE 2980
jedi@server:~$ _
```

5

# poll()

- Wait for one or more file descriptors to become ready for use
- Positives
    - POSIX (Portable Operating System Interface, from IEEE)
        - can be used on Linux, BSD flavors, *etc.*
- Negatives
    - Does not scale to many file descriptors

```c
#include <stdio.h>
#include <unistd.h>
#include <sys/poll.h>


#define TIMEOUT 5


int main (void)
{
        struct pollfd fds[2];
        int ret;

        /* watch stdin for input */
        fds[0].fd = STDIN_FILENO;
        fds[0].events = POLLIN;

        /* watch stdout for ability to write */
        fds[1].fd = STDOUT_FILENO;
        fds[1].events = POLLOUT;

        ret = poll(fds, 2, TIMEOUT * 1000);
```

7

# ppoll()

- A race condition can occur if there are any signal handlers registered, ppoll() atomically handles signals, applies a sigmask, and saves new incoming signals to the end

- More details below in description of pselect()

# select()

- Like poll(), but older and clunkier

```
      int select(int nfds, fd_set *readfds, fd_set *writefds,
                 fd_set *exceptfds, struct timeval *timeout);

      void FD_CLR(int fd, fd_set *set);
      int  FD_ISSET(int fd, fd_set *set);
      void FD_SET(int fd, fd_set *set);
      void FD_ZERO(fd_set *set);

      int pselect(int nfds, fd_set *readfds, fd_set *writefds,
                  fd_set *exceptfds, const struct timespec *timeout,
                  const sigset_t *sigmask);

   Feature Test Macro Requirements for glibc (see feature_test_macros(7)):

      pselect(): _POSIX_C_SOURCE >= 200112L

DESCRIPTION
      select() allows a program to monitor multiple file descriptors, waiting
      until one or more of the file descriptors become "ready" for some class
      of I/O operation (e.g., input possible).  A file descriptor is  consid-
      ered  ready  if it is possible to perform a corresponding I/O operation
      (e.g., read(2), or a sufficiently small write(2)) without blocking.

      select() can monitor only file descriptors numbers that are  less  than
      FD_SETSIZE;  poll(2)  and  epoll(7)  do  not have this limitation.  See
      BUGS.
```
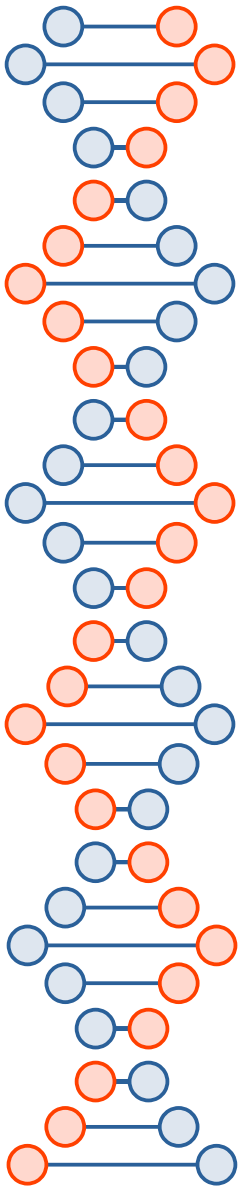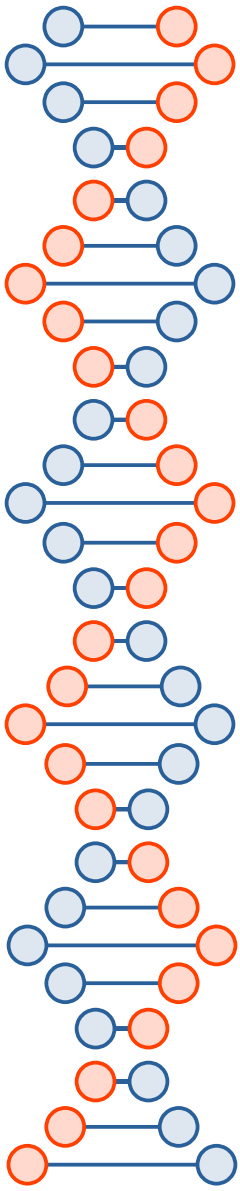
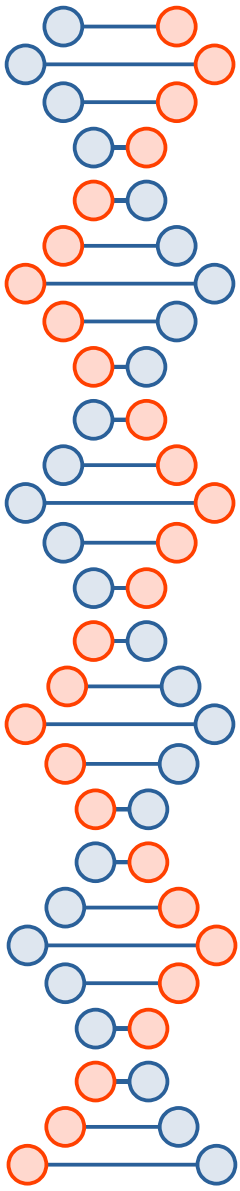Manual page select(2) line 9 (press h for help or q to quit)

**pselect()**

The **pselect**() system call allows an application to safely wait until either a file descriptor becomes ready or until a signal is caught.

The operation of **select**() and **pselect**() is identical, other than these three differences:

- **select**() uses a timeout that is a <u>struct</u> <u>timeval</u> (with seconds and microseconds), while **pselect**() uses a <u>struct</u> <u>timespec</u> (with seconds and nanoseconds).

- **select**() may update the <u>timeout</u> argument to indicate how much time was left. **pselect**() does not change this argument.

- **select**() has no <u>sigmask</u> argument, and behaves as **pselect**() called with NULL <u>sigmask</u>.

<u>sigmask</u> is a pointer to a signal mask (see **sigprocmask**(2)); if it is not NULL, then **pselect**() first replaces the current signal mask by the one pointed to by <u>sigmask</u>, then does the "select" function, and then restores the original signal mask. (If <u>sigmask</u> is NULL, the signal mask is not modified during the **pselect**() call.)

Other than the difference in the precision of the <u>timeout</u> argument, the following **pselect**() call:

    ready = pselect(nfds, &readfds, &writefds, &exceptfds,

Manual page select(2) line 132 (press h for help or q to quit)

1

Other than the difference in the precision of the <u>timeout</u> argument, the following **pselect**() call:

```
ready = pselect(nfds, &readfds, &writefds, &exceptfds,
                timeout, &sigmask);
```

is equivalent to <u>atomically</u> executing the following calls:

```
sigset_t origmask;

pthread_sigmask(SIG_SETMASK, &sigmask, &origmask);
ready = select(nfds, &readfds, &writefds, &exceptfds, timeout);
pthread_sigmask(SIG_SETMASK, &origmask, NULL);
```

The reason that **pselect**() is needed is that if one wants to wait for either a signal or for a file descriptor to become ready, then an atomic test is needed to prevent race conditions. (Suppose the signal handler sets a global flag and returns. Then a test of this global flag followed by a call of **select**() could hang indefinitely if the signal arrived just after the test but just before the call. By contrast, **pselect**() allows one to first block signals, handle the signals that have come in, then call **pselect**() with the desired <u>sigmask</u>, avoiding the race.)

**The timeout**
        The <u>timeout</u> argument for **select**() is a structure of the following type:

2

# epoll()

- Negatives
  - Not POSIX, Linux-specific
  - Slightly more complex to use than poll()
- Positives

| Number of File Descriptors | poll() CPU time | select() CPU time | epoll() CPU time |
|---|---|---|---|
| 10 | 0.61 | 0.73 | 0.41 |
| 100 | 2.9 | 3 | 0.42 |
| 1000 | 35 | 35 | 0.53 |
| 10000 | 990 | 930 | 0.66 |

The Linux Programming Interface, section 63.4.5
https://suchprogramming.com/epoll-in-3-easy-steps/

Examples from

https://suchprogramming.com/epoll-in-3-easy-steps/

...

# Step 1: Create epoll file descriptor

First I'll go through the process of just creating and closing an epoll instance.

```c
#include <stdio.h>      // for fprintf()
#include <unistd.h>     // for close()
#include <sys/epoll.h>  // for epoll_create1()

int main()
{
        int epoll_fd = epoll_create1(0);

        if (epoll_fd == -1) {
                fprintf(stderr, "Failed to create epoll file descriptor\n");
                return 1;
        }

        if (close(epoll_fd)) {
                fprintf(stderr, "Failed to close epoll file descriptor\n");
                return 1;
        }

        return 0;
}
```

## Step 2: Add file descriptors for epoll to watch

The next thing to do is tell epoll what file descriptors to watch and what kinds of events to watch for. In this example I'll use one of my favorite file descriptors in Linux, good ol' file descriptor `0` (also known as Standard Input).

```c
#include <stdio.h>      // for fprintf()
#include <unistd.h>     // for close()
#include <sys/epoll.h>  // for epoll_create1(), epoll_ctl(), struct epoll_event

int main()
{
        struct epoll_event event;
        int epoll_fd = epoll_create1(0);

        if (epoll_fd == -1) {
                fprintf(stderr, "Failed to create epoll file descriptor\n");
                return 1;
        }

        event.events = EPOLLIN;
        event.data.fd = 0;

        if (epoll_ctl(epoll_fd, EPOLL_CTL_ADD, 0, &event)) {
                fprintf(stderr, "Failed to add file descriptor to epoll\n");
                close(epoll_fd);
                return 1;
        }

        if (close(epoll_fd)) {
                fprintf(stderr, "Failed to close epoll file descriptor\n");
                return 1;
        }
        return 0;

}
```
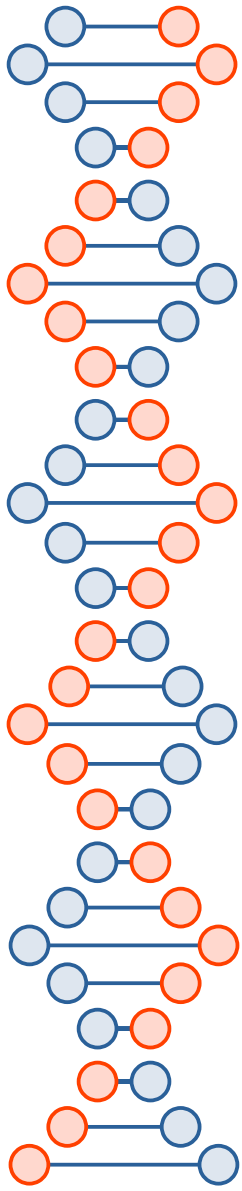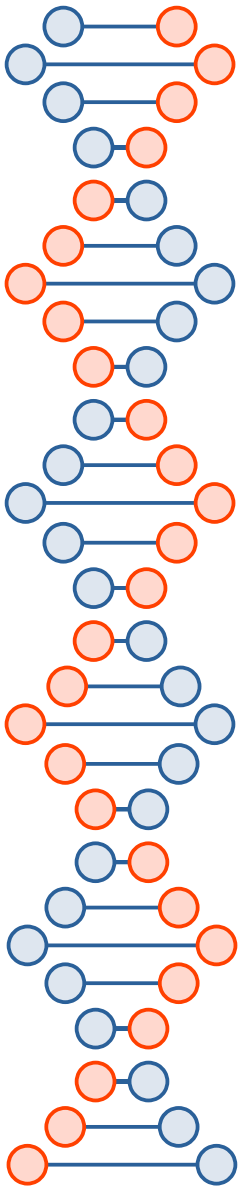
## Step 3: Profit

That's right! We're almost there. Now let epoll do it's magic.

```c
while (running) {
    printf("\nPolling for input...\n");
    event_count = epoll_wait(epoll_fd, events, MAX_EVENTS, 30000);
    printf("%d ready events\n", event_count);
    for (i = 0; i < event_count; i++) {
        printf("Reading file descriptor '%d' -- ", events[i].data.fd
        bytes_read = read(events[i].data.fd, read_buffer, READ_SIZE)
        printf("%zd bytes read.\n", bytes_read);
        read_buffer[bytes_read] = '\0';
        printf("Read '%s'\n", read_buffer);

        if(!strncmp(read_buffer, "stop\n", 5))
            running = 0;
    }
}
```

Why is epoll() faster?

https://en.wikipedia.org/wiki/Red%E2%80%93black_tree

## Red–black tree

| | |
|---|---|
| **Type** | Tree |
| **Invented** | 1978 |
| **Invented by** | Leonidas J. Guibas and Robert Sedgewick |

### Complexities in big O notation

**Space complexity**

| **Space** | $O(n)$ |
|---|---|

**Time complexity**

| **Function** | **Amortized** | **Worst Case** |
|---|---|---|
| **Search** | $O(\log n)$[1] | $O(\log n)$[1] |
| **Insert** | $O(1)$[2] | $O(\log n)$[1] |
| **Delete** | $O(1)$[2] | $O(\log n)$[1] |



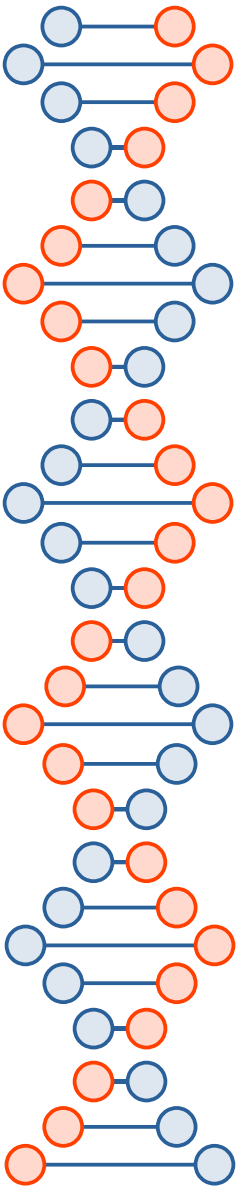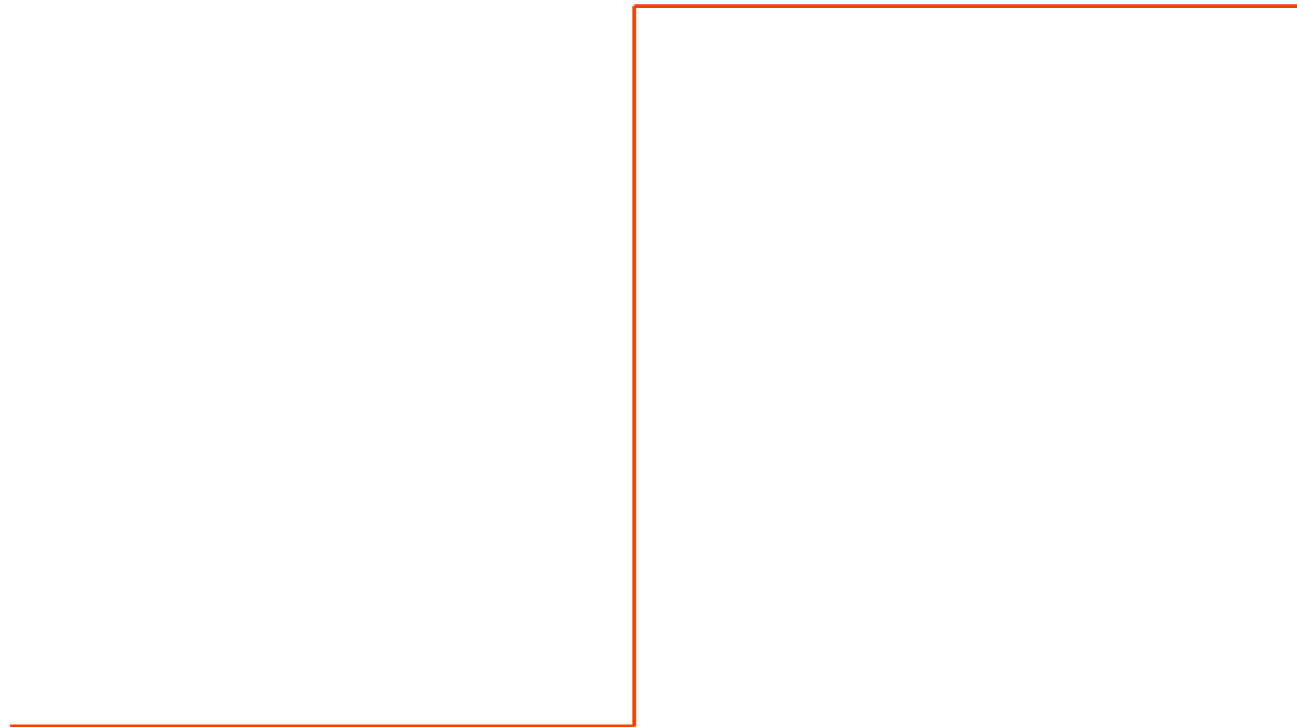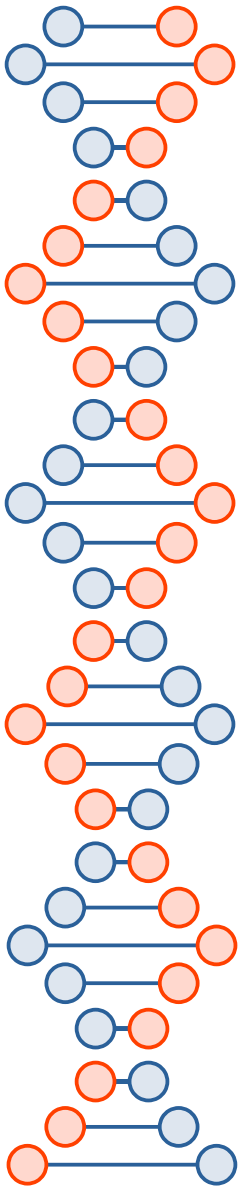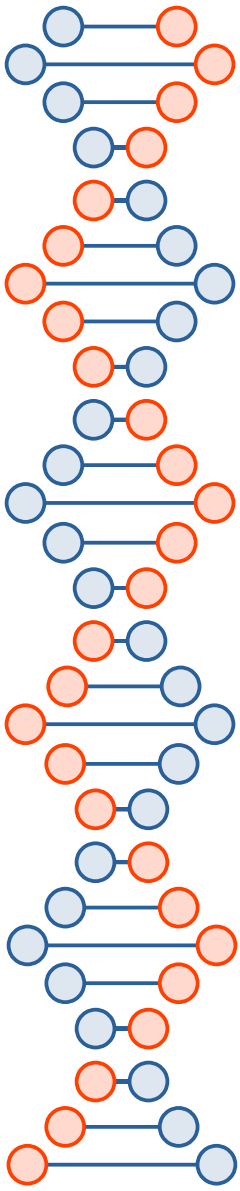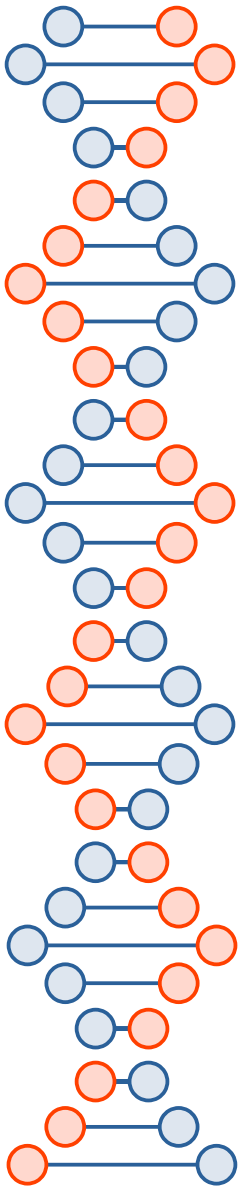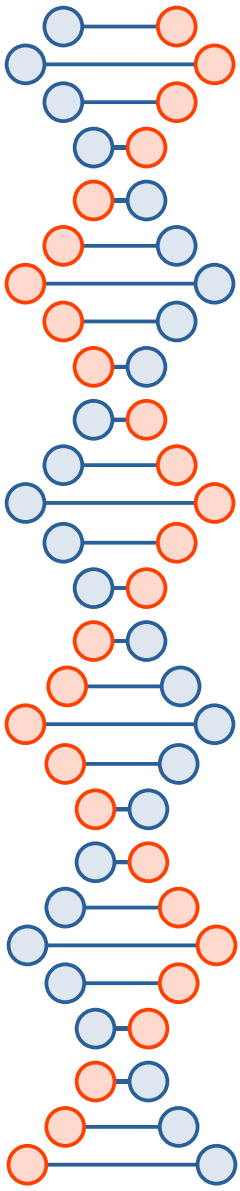Figure 1: ... with explicit NIL leaves

19

Edge- *vs.* level-triggered?

## Level-triggered and edge-triggered

The **epoll** event distribution interface is able to behave both as edge-triggered (ET) and as level-triggered (LT). The difference between the two mechanisms can be described as follows. Suppose that this scenario happens:

1. The file descriptor that represents the read side of a pipe (<u>rfd</u>) is registered on the **epoll** instance.

2. A pipe writer writes 2 kB of data on the write side of the pipe.

3. A call to **epoll_wait**(2) is done that will return <u>rfd</u> as a ready file descriptor.

4. The pipe reader reads 1 kB of data from <u>rfd</u>.

5. A call to **epoll_wait**(2) is done.

If the <u>rfd</u> file descriptor has been added to the **epoll** interface using the **EPOLLET** (edge-triggered) flag, the call to **epoll_wait**(2) done in step **5** will probably hang despite the available data still present in the file input buffer; meanwhile the remote peer might be expecting a response based on the data it already sent. The reason for this is that edge-triggered mode delivers events only when changes occur on the monitored file descriptor. So, in step **5** the caller might end up waiting for some data that is already present inside the input buffer. In the above example, an event on <u>rfd</u> will be generated because of the write done in **2** and the event is consumed in **3**. Since the read operation done in **4** does not consume the whole buffer data,

Manual page epoll(7) line 42 (press h for help or q to quit)

2

5. A call to **epoll_wait**(2) is done.

If the _rfd_ file descriptor has been added to the **epoll** interface using the **EPOLLET** (edge-triggered) flag, the call to **epoll_wait**(2) done in step **5** will probably hang despite the available data still present in the file input buffer; meanwhile the remote peer might be expecting a response based on the data it already sent. The reason for this is that edge-triggered mode delivers events only when changes occur on the monitored file descriptor. So, in step **5** the caller might end up waiting for some data that is already present inside the input buffer. In the above example, an event on _rfd_ will be generated because of the write done in **2** and the event is consumed in **3**. Since the read operation done in **4** does not consume the whole buffer data, the call to **epoll_wait**(2) done in step **5** might block indefinitely.

An application that employs the **EPOLLET** flag should use nonblocking file descriptors to avoid having a blocking read or write starve a task that is handling multiple file descriptors. The suggested way to use **epoll** as an edge-triggered (**EPOLLET**) interface is as follows:

a) with nonblocking file descriptors; and

b) by waiting for an event only after **read**(2) or **write**(2) return **EAGAIN**.

By contrast, when used as a level-triggered interface (the default, when **EPOLLET** is not specified), **epoll** is simply a faster **poll**(2), and can be used wherever the latter is used since it shares the same semantics.

Manual page epoll(7) line 57 (press h for help or q to quit)

23

Since even with edge-triggered **epoll**, multiple events can be generated upon receipt of multiple chunks of data, the caller has the option to specify the **EPOLLONESHOT** flag, to tell **epoll** to disable the associated file descriptor after the receipt of an event with **epoll_wait**(2). When the **EPOLLONESHOT** flag is specified, it is the caller's responsibility to rearm the file descriptor using **epoll_ctl**(2) with **EPOLL_CTL_MOD**.

If multiple threads (or processes, if child processes have inherited the **epoll** file descriptor across **fork**(2)) are blocked in **epoll_wait**(2) waiting on the same epoll file descriptor and a file descriptor in the interest list that is marked for edge-triggered (**EPOLLET**) notification becomes ready, just one of the threads (or processes) is awoken from **epoll_wait**(2). This provides a useful optimization for avoiding "thundering herd" wake-ups in some scenarios.

**Interaction with autosleep**
    If the system is in **autosleep** mode via _/sys/power/autosleep_ and an event happens which wakes the device from sleep, the device driver will keep the device awake only until that event is queued. To keep the device awake until the event has been processed, it is necessary to use the **epoll_ctl**(2) **EPOLLWAKEUP** flag.

    When the **EPOLLWAKEUP** flag is set in the **events** field for a _struct epoll_event_, the system will be kept awake from the moment the event is queued, through the **epoll_wait**(2) call which returns the event until the subsequent **epoll_wait**(2) call. If the event should keep the system awake

Manual page epoll(7) line 83 (press h for help or q to quit)

4

- "In computer science, the thundering herd problem occurs when a large number of processes or threads waiting for an event are awoken when that event occurs, but only one process is able to handle the event. When the processes wake up, they will each try to handle the event, but only one will win. All processes will compete for resources, possibly freezing the computer, until the herd is calmed down again."

25

# Use case: Tor

- Overlay network that provides anonymity and censorship resistance

- Easy to use (Tor browser), open source, friendly to academics (and lots of data)

- Uses epoll() for Linux and kqueue() for BSD flavors

How Tor Works: 2

Legend:
- Tor node
- unencrypted link
- encrypted link

Alice

Step 2: Alice's Tor client picks a random path to destination server. Green links are encrypted, red links are in the clear.

Jane

Dave

Bob

https://en.wikipedia.org/wiki/Tor_(network)

27

Some metrics from
https://metrics.torproject.org/

…

Directly connecting users

The Tor Project - https://metrics.torproject.org/

29

# Number of relays



Relays
Bridges

The Tor Project - https://metrics.torproject.org/

# Advertised and consumed bandwidth by relay flags



The Tor Project - https://metrics.torproject.org/

31

# What is kqueue()

- Similar to epoll(), but for BSD flavors
- "Kqueue allows one to batch modify watcher states and to retrieve watcher states in a single system call. With epoll, you have to call a system call for every modification. Kqueue also allows one to watch for things like filesystem changes and process state changes, epoll is limited to socket/pipe I/O only."

    --asomiv, https://news.ycombinator.com/item?id=3028687

- Linux has inotify()
- libuv and libevent support kqueue(), epoll(), and alternatives such as Solaris I/O completion ports, Windows IOCP, *etc.*

We've mentioned Solaris a few times, now seeing a difference from Linux *vs*. BSD, and we've largely ignored Windows this semester. Now's a good time for a 5-minute digression into OS history...

1969
1971 to 1973
1974 to 1975
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001 to 2004
2005
2006 to 2007
2008
2009
2010
2011
2012 to 2015
2016
2017
2018
2019

1969
1971 to 1973
1974 to 1975
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001 to 2004
2005
2006 to 2007
2008
2009
2010
2011
2012 to 2015
2016
2017
2018
2019

Unnamed PDP-7 operating system

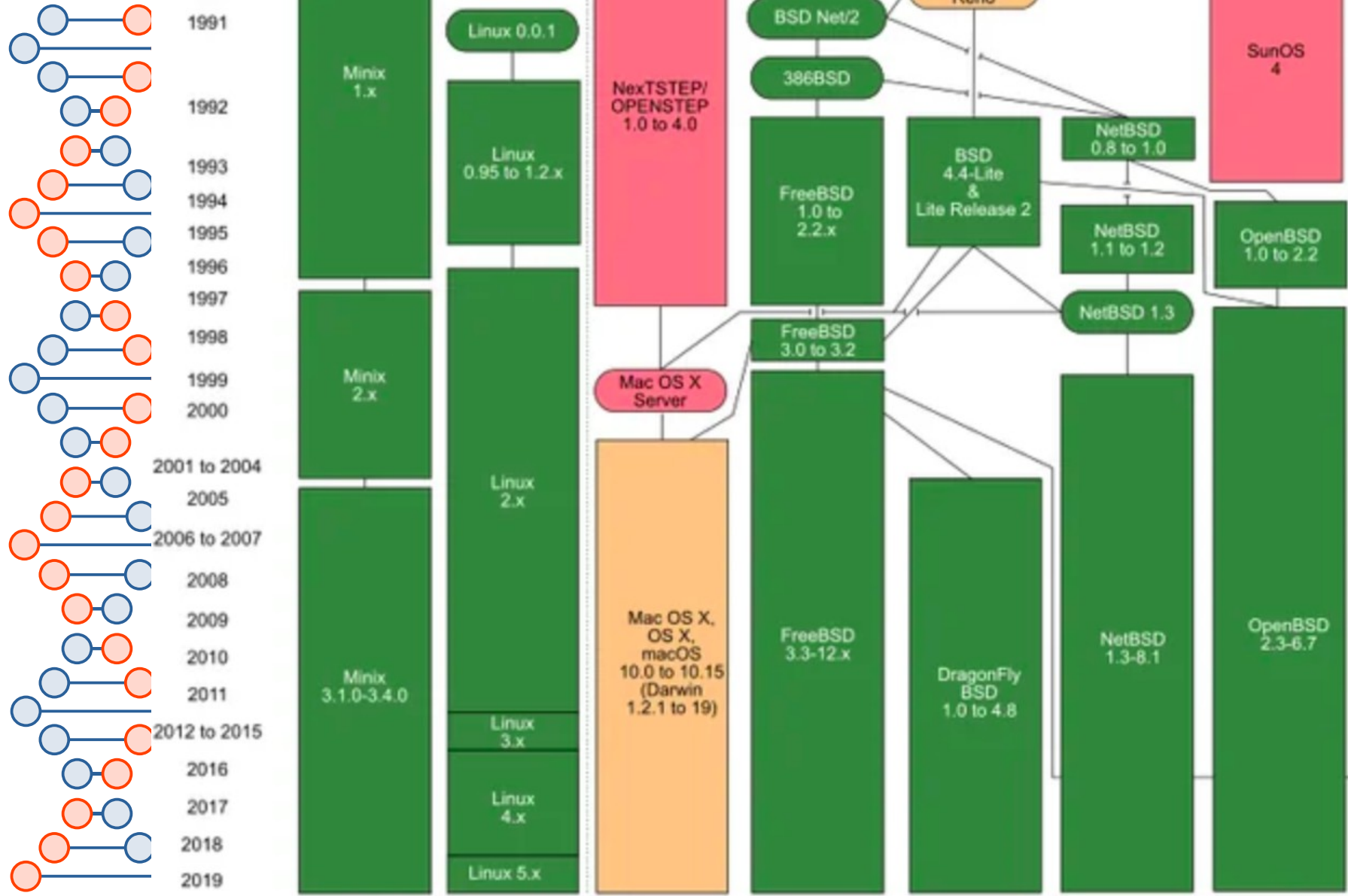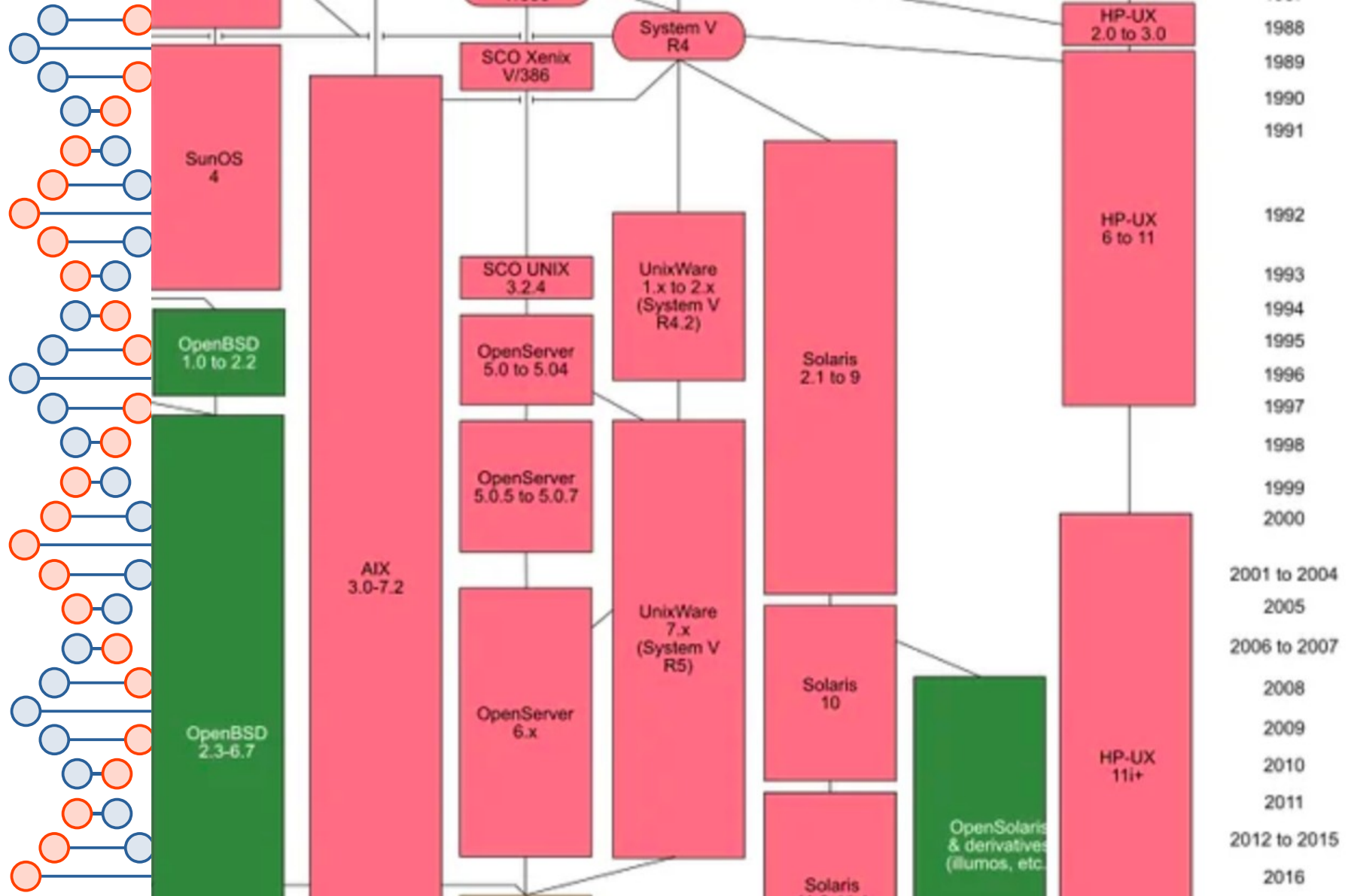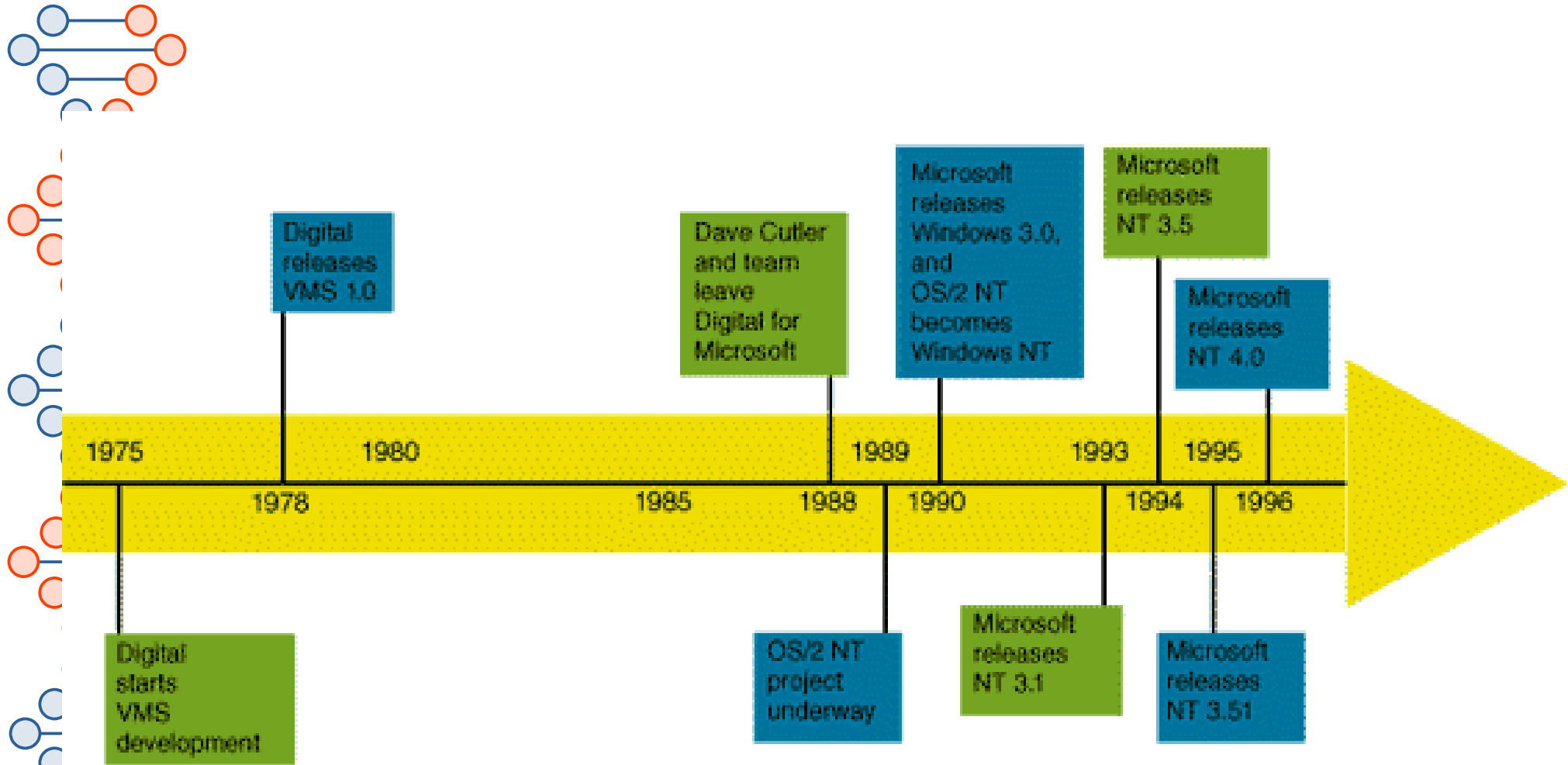Unix Version 1 to 4

Unix Version 5 to 6

PWB/Unix

BSD 1.0 to 2.0

Unix Version 7

Unix/32V

https://www.reddit.com/r/linux/comments/huhqrh/unix_family_tree/

**Open source**
**Mixed/shared source**
**Closed source**

BSD 3.0 to 4.1

Xenix 1.0 to 2.3

System III

BSD 4.2

SunOS 1 to 1.1

Xenix 3.0

System V R1 to R2

Unix Version 8

SCO Xenix

Unix-like systems

BSD 4.3

SunOS 1.2 to 3.0

AIX 1.0

SCO Xenix V/286

System V R3

HP-UX 1.0 to 1.2

Unix 9 and 10 (last versions from Bell Labs)

BSD 4.3 Tahoe

SCO Xenix V/386

HP-UX 2.0 to 3.0

System V R4

BSD 4.3 Reno

BSD Net/1

SCO Xenix V/386

Minix 1.x

Linux 0.0.1

BSD Net/2

SunOS 4

NexTSTEP/ OPENSTEP 1.0 to 4.0

386BSD

Linux 0.95 to 1.2.x

SCO UNIX 3.2.4

UnixWare 1.x to 2.x (System V R4.2)

HP-UX 6 to 11

FreeBSD 1.0 to 2.2.x

BSD 4.4-Lite & Lite Release 2

NetBSD 0.8 to 1.0

Minix 2.x

Mac OS X Server

NetBSD 1.1 to 1.2

OpenBSD 1.0 to 2.2

OpenServer 5.0 to 5.04

Solaris 2.1 to 9

Linux 2.x

FreeBSD 3.0 to 3.2

NetBSD 1.3

OpenServer 5.0.5 to 5.0.7

AIX 3.0-7.2

Mac OS X, OS X, macOS 10.0 to 10.15 (Darwin 1.2.1 to 19)

UnixWare 7.x (System V R5)

Minix 3.1.0-3.4.0

FreeBSD 3.3-12.x

DragonFly BSD 1.0 to 4.8

NetBSD 1.3-8.1

OpenBSD 2.3-6.7

OpenServer 6.x

Solaris 10

Linux 3.x

Solaris 11.0-11.4

OpenSolaris & derivatives (illumos, etc)

HP-UX 11i+

Linux 4.x

OpenServer 10.x

34

Linux 5.x

36

SunOS 4

OpenBSD 1.0 to 2.2

OpenBSD 2.3-6.7

AIX 3.0-7.2

SCO Xenix V/386

SCO UNIX 3.2.4

OpenServer 5.0 to 5.04

OpenServer 5.0.5 to 5.0.7

OpenServer 6.x

System V R4

UnixWare 1.x to 2.x (System V R4.2)

UnixWare 7.x (System V R5)

Solaris 2.1 to 9

Solaris 10

OpenSolaris & derivatives (illumos, etc.)

Solaris

HP-UX 2.0 to 3.0

HP-UX 6 to 11

HP-UX 11i+

1988
1989
1990
1991

1992

1993
1994
1995
1996
1997
1998
1999
2000
2001 to 2004
2005
2006 to 2007
2008
2009
2010
2011
2012 to 2015
2016

37

Timeline:

- 1975 — Digital starts VMS development
- 1978 — Digital releases VMS 1.0
- 1980
- 1985 — Dave Cutler and team leave Digital for Microsoft
- 1988 — OS/2 NT project underway
- 1989 — Microsoft releases Windows 3.0, and OS/2 NT becomes Windows NT
- 1990 — Microsoft releases NT 3.1
- 1993 — Microsoft releases NT 3.5
- 1994 — Microsoft releases NT 3.51
- 1995 — Microsoft releases NT 4.0
- 1996

https://www.tech-insider.org/windows/research/1998/images/figure_01.gif

38

# CPM - 1974

https://en.wikipedia.org/wiki/CP/M#/media/File:CPM-86.png

February 3, 1976

## An Open Letter to Hobbyists

To me, the most critical thing in the hobby market right now is the lack of good software courses, books and software itself. Without good software and an owner who understands programming, a hobby computer is wasted. Will quality software be written for the hobby market?

Almost a year ago, Paul Allen and myself, expecting the hobby market to expand, hired Monte Davidoff and developed Altair BASIC. Though the initial work took only two months, the three of us have spent most of the last year documenting, improving and adding features to BASIC. Now we have 4K, 8K, EXTENDED, ROM and DISK BASIC. The value of the computer time we have used exceeds $40,000.

[See the fulll letter at https://en.wikipedia.org/wiki/File:Bill_Gates_Letter_to_Hobbyists_ocr.pdf]
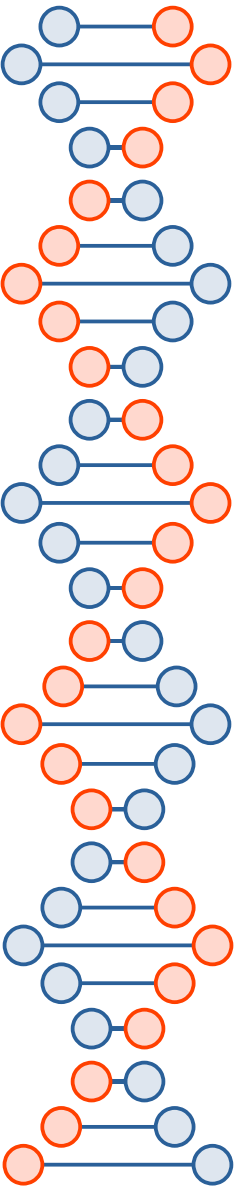
I would appreciate letters from any one who wants to pay up, or has a suggestion or comment. Just write me at 1180 Alvarado SE, #114, Albuquerque, New Mexico, 87108. Nothing would please me more than being able to hire ten programmers and deluge the hobby market with good software.

*Bill Gates*

Bill Gates
General Partner, Micro-Soft

40

# QDOS - 1979

```
A:asm mon

Seattle Computer Products
Copyright 1979,80,81 by Se



Error Count =      0

A:hex2bin mon

A:_
```

41

# MS-DOS - 1981

```
Starting MS-DOS...

C:\>_
```