API. 海代码层的接口。 ABI 二进制指心· 通过具证、编书器、被持器、

两

POSIX 和SUS 翻绕机接口粉的 UMIX

C 25 AMS 2 60 250C-

gcc c语言 gtt Ctt语言

LSB Linux towards

一切简之件

fd 之件描述符 int

文件条件为

支持多个进行工艺的一个是人没有限到文件的专家的的

要到花. mode 洲植野 专之体美貌

国录、可选加标的 mode 的政射·

链接加水、

国地南部松 V

对国民的操作: 安加放转, 柳阳的转

超越接,多台名特映新到一个inode.

符号发射等: 快拍为关. (超频前的的made 和数据)

由 扫描全能王 扫描创建

件部外之件 (以之件表主的内本文对象) 快没备好 作为进数组为同人硬盘,等 随至为的 子符没有这件. 铁性新到去openhand BUZ. 布名等道 7-270/台ZPC (在文件) 普通普遍在的存在 []NIX核套接谷 使用Socher 2件交接

文件的物分为空间·(统一) 七层加,翻陷、链数细渐新数 挂我气 挂我到多知识的特定位置 历起(分级多行文件分配) · sector Box. block to

H43

ELF 哥姆的哥拉特之件(习这份的代码) 效极. 数据段 b4542.

线程有有格局的是量的栈

Pld

进程权 根的法 的社为各场份 进程企业并不被删除,等待交进经查和依法大删除 用户和组

227 721月日12月美族 (英定 n2d)

Yout vid = 0 effective vid saved vid - J's vid

每个那层于一支多个组: 920

权限

文件权限

外福烟的校假凌,百,批约) 雄似的积 三千子经、共9亿、保存在incele中 YWX 者俱地

傷

平6年多通知机到. 数值率量十之格

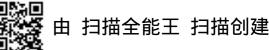
IPC.

管道, 都管道, 作客量, 净象队列共和待, 快到那些问至分

进件

裕没处理.

一. 特殊空星 errnor ta < errnu. L)+ exten int errno!



Zstolio h> + perrorl). straror(). straror\_r(). 2/19 1/0 效件表。(参). fd. int Cookes 打开这件返回fd: 进程 至了了一fd 0; stdin 1: Stdant. 2:stderr 海野到: berboard. —screen. 理之句. fol 可治河 没名之件, 营道, fustesses 7-170 Socket 打入文件 open ( name, flag, mode) read ( )

0-1-10

标准3/0 是C标准每提供的一个用户级冲存

41. 分数/聚集 1/0

一次润用对新线冲压输入输出 引新礼 (支持的性) 可向量 2/0 (共产之战性2/0) reado() wider()

iovec 结构体 描述一个线冲区,(我)一组段特何是 Count 较大:内核动态分配发析 · 较小:在松上创建个级数组·

4.2. Event Pell. (Linux 查有).
Pull 和solect 李要证的 fol 根例 性能差.

Frent pul 三: 1. 新知如何别, 2, 深如倒得fol, s. 事件错.

1-epoll-crewtell) 这四fd. 监视农场里用close()美洲

2. apoll-att() tox Att) Forfol. exoll-event 45 take.

3. epoll\_waste). (timeont

沙锅(新发,各件的校。

边缘就发。每两些各件的操作结束后返回。

每件条件:户满足多件的操作一发生的效图(黑人水)

45、有治映射·(进程地址空间) 将之件印料到内存中的到内部,直接通过内分为约约 mmap()



元是内存管理单元/MMU)的\*年基单位, Long page\_Size = Sysconf(\_SC-PAGESLZE); \*\*\*\*\* <asmpages.h) PAGIE\_STZT:

信号? format

munphape) Aug mmapelity 47.

mmap() 不太恭来

mremap() 词程映射大

油运mrengp,mmap等实现内部部。mprotect 改变物河校限

通过映新了的文件.

msyrc. (李河fsync) 将映射哥回磁金. madvise

放废、满煤块办会,提前海和了一块

44. 省殖文件110程方.
Posix-fadvise() for readahard()

4.5. Synchronous 海岛操作线束后这回。 asy-Synchronoized·海岛排作用排料显微型。nonsy...
Law. Lo 虽为 2/0.

4.6. 1/0词及器和 2/0性能 Seek, 石狐从移动 Sems以上 1/0词及器, 油运辑 1/0清新的顺台和时间点, 担制比较 (1)成为 Seek和的)

拉面 3每年址: 柱面,3每头,扇丘,三元组合为块色(1BA) 3如道 鹰片

2月的强度. (物理状的整数法)

新 多种物的 40 满花的个

英满水,(分次返回最新数据). 基数据不在双线存中,从磁整线出本阻塞(deby) 避免排作若来满地激动

Co. Deadline 210 河流

三千队列 核焰队列 淀刊和 写刊FO 排为 Sooms SS

②·Anticipatory 1/0 调度. (期望. 对凌橇 响应公敦后等待的与期望连续的淡

③中日 110 (完全从平限到)

新进程和队列和附明的批准

(H) Noop 2/0 三新,不排序.

SSD 用 Map 或 CFQ 优化 2/0性能 用空间 210 洞餐

2种的技术的20满年进份排序程务。 1抽路经排序。 2抽油超标符 3抽物组体排作 Solder Stelerr