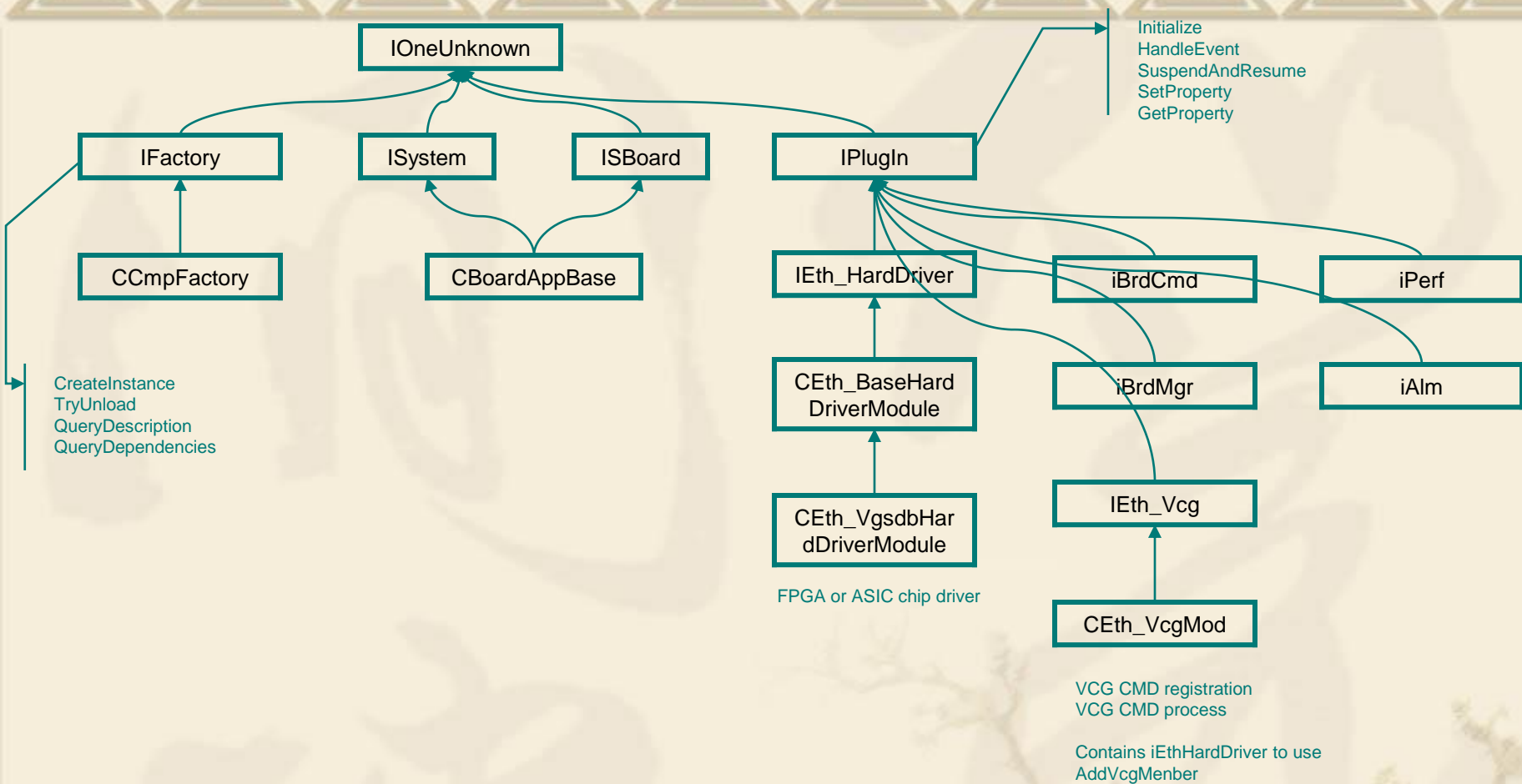1. A COM application (Harbour Networks)
2. An Network OS (Huawei VRP)
3. C++ advanced features (OOP, C++11)
4. SOCKET&TCP illustration (NetBSD)
5. OS concepts

Harold Mei

IOneUnknown

IFactory

ISystem

ISBoard

IPlugIn

Initialize
HandleEvent
SuspendAndResume
SetProperty
GetProperty

CCmpFactory

CBoardAppBase

IEth_HardDriver

iBrdCmd

iPerf

CreateInstance
TryUnload
QueryDescription
QueryDependencies

CEth_BaseHard
DriverModule

iBrdMgr

iAlm

CEth_VgsdbHar
dDriverModule

IEth_Vcg

FPGA or ASIC chip driver

CEth_VcgMod

VCG CMD registration
VCG CMD process

Contains iEthHardDriver to use
AddVcgMenber

| | |
|---|---|
| Core: IP/MPLS/VPN/RM | SCP |
| | |
| IFNET | IPV4 FWD | SMP |
| PPP/ATM/ETH | | CFM |
| LinkLayer | Security IPSec/ACL/FW/QOS | IC |
| | | SRM |

SSP: Mem/Timer/RPC/IPC/Process/Load/OS

## 1. Encapsulation:

public, private, protected;
All are inherited, but the accessibilities for those three are different.

## 2. Inheritance:

Overload: same function name, different parameters. Could be class members or independent functions.
Overwrite: virtual member functions, **polymorphism & RTTI**.
Override: for non-virtual member functions.

Virtual member function: can be called by a base class object pointer, if the pointer refers to objects of a class which the function is implemented.
Non-virtual member function: can not be called by a base class object.

## 3. polymorphism

A base pointer refers to a descendent, hence when similar logic process (generalization) need to call a virtual function from the base class, it will call the **nearest** overwritten virtual function from the exact object class.

## 4. static member

Even though static member is defined in class, it is independent of any objects, and it should be initialized outside the constructor.

## 5. Constructor & Destructor

Copy constructor:
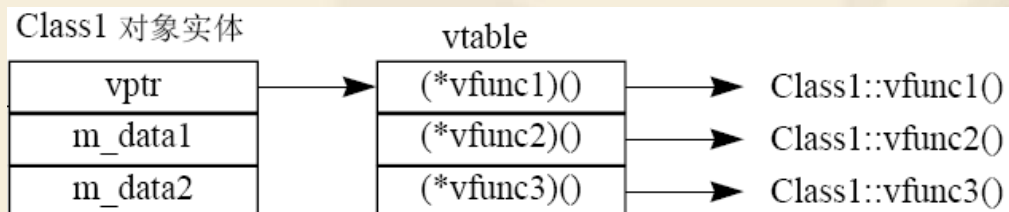Assignment constructor: avoid default constructors.
Operator xx():
Implicitly and explicitly called in programs.

## 6. Exception handling

## 7. Template

## *1. Memory type
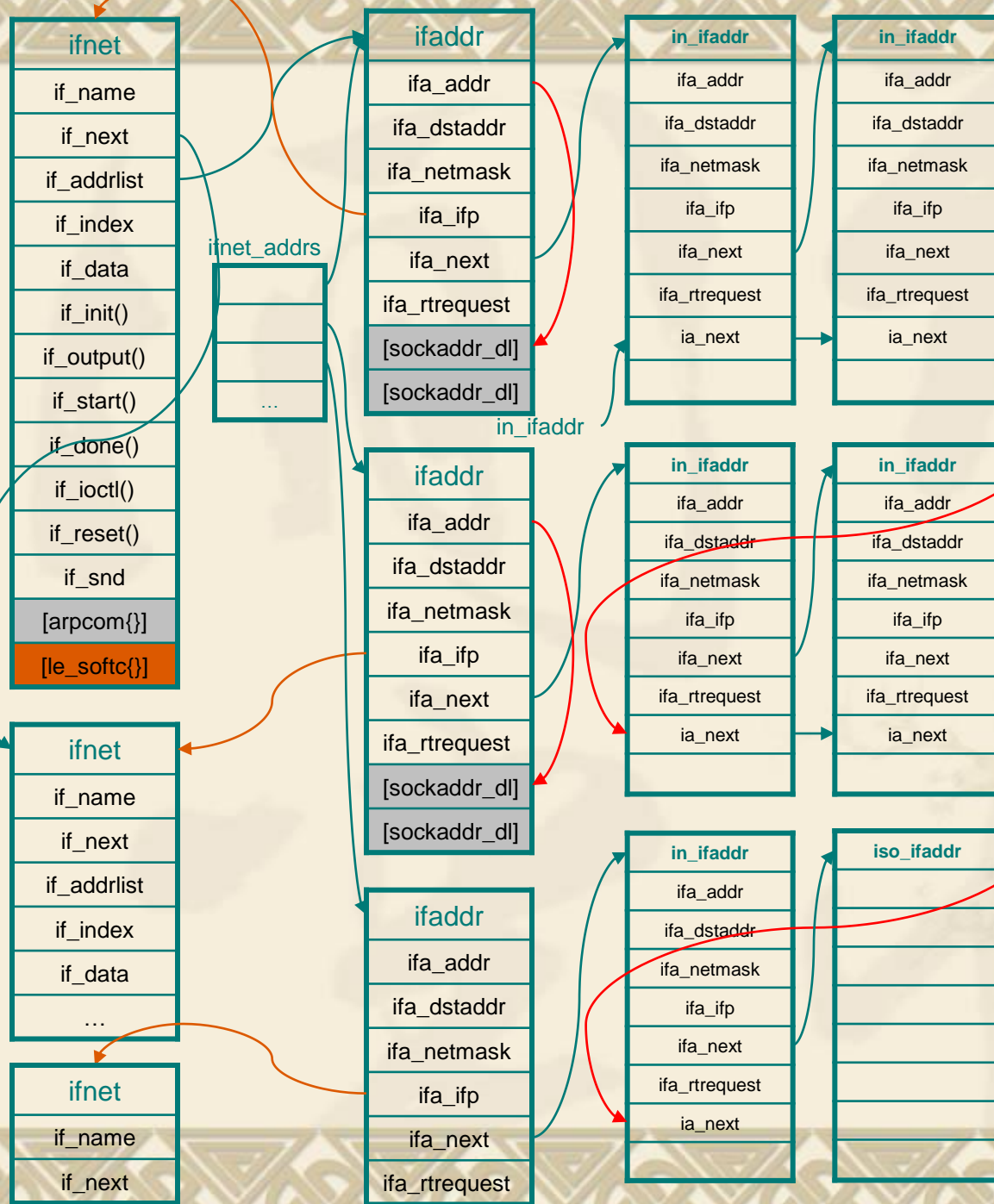
Static memory; stack; heap.
volatile;

**1、泛化(继承)**
**2、实现**
3.1单向关联
3.2双向关联
3.3聚合关系
3.4组合关系
**4、依赖**

**ifnet**

| |
|---|
| if_name |
| if_next |
| if_addrlist |
| if_index |
| if_data |
| if_init() |
| if_output() |
| if_start() |
| if_done() |
| if_ioctl() |
| if_reset() |
| if_snd |
| [arpcom{}] |
| [le_softc{}] |

ifnet_addrs

| |
|---|
| |
| |
| |
| ... |

**ifaddr**

| |
|---|
| ifa_addr |
| ifa_dstaddr |
| ifa_netmask |
| ifa_ifp |
| ifa_next |
| ifa_rtrequest |
| [sockaddr_dl] |
| [sockaddr_dl] |

in_ifaddr

**in_ifaddr**

| |
|---|
| ifa_addr |
| ifa_dstaddr |
| ifa_netmask |
| ifa_ifp |
| ifa_next |
| ifa_rtrequest |
| ia_next |

**in_ifaddr**

| |
|---|
| ifa_addr |
| ifa_dstaddr |
| ifa_netmask |
| ifa_ifp |
| ifa_next |
| ifa_rtrequest |
| ia_next |

**ifaddr**

| |
|---|
| ifa_addr |
| ifa_dstaddr |
| ifa_netmask |
| ifa_ifp |
| ifa_next |
| ifa_rtrequest |
| [sockaddr_dl] |
| [sockaddr_dl] |

**in_ifaddr**

| |
|---|
| ifa_addr |
| ifa_dstaddr |
| ifa_netmask |
| ifa_ifp |
| ifa_next |
| ifa_rtrequest |
| ia_next |

**in_ifaddr**

| |
|---|
| ifa_addr |
| ifa_dstaddr |
| ifa_netmask |
| ifa_ifp |
| ifa_next |
| ifa_rtrequest |
| ia_next |

**ifnet**

| |
|---|
| if_name |
| if_next |
| if_addrlist |
| if_index |
| if_data |
| … |

**ifaddr**

| |
|---|
| ifa_addr |
| ifa_dstaddr |
| ifa_netmask |
| ifa_ifp |
| ifa_next |
| ifa_rtrequest |

**in_ifaddr**

| |
|---|
| ifa_addr |
| ifa_dstaddr |
| ifa_netmask |
| ifa_ifp |
| ifa_next |
| ifa_rtrequest |
| ia_next |

**iso_ifaddr**

| |
|---|
| |
| |
| |
| |
| |
| |
| |

**ifnet**

| |
|---|
| if_name |
| if_next |

```
main(framep)
void *framep;
{
cpu_startup();
…
for (pdev = pdevinit; pdev->pdev_attach !=
NULL; pdev++)
    (*pdev->pdev_attach)(pdev->pdev_count);
//if_attach(); ether_attach();
…
ifinit();
domaininit();
scheduler();
}
```

IOCTL命令:
0, SIOCGIFCONF: 获得接口配置
1, SIOCGIFFLAGS:获得接口标志
2, SIOCGIFMETRIC: 获得接口度量
3, SIOCSIFFLAGS: 设置接口标志(LINK status conf)
4, SIOCSIFMETRIC: 设置接口度量
5, SIOCG(S)IFADDR
6, SIOCG(S)IFNETMASK
7, SIOCG(S)IFDSTADDR
8, SIOCG(S)IFBRDADDR
9, SIOCDIFADDR
10, SIOCAIFADDR

| **sockaddr_dl** | **sockaddr** | **sockaddr_in** |
|---|---|---|
| sdl_len | len | len |
| sdl_family | family | family |
| sdl_index | data | port(2 bytes) |
| sdl_type | (14 bytes) | addr |
| sdl_nlen | | (4 bytes) |
| sdl_alen | | zero |
| sdl_slen | | (8 bytes) |
| sdl_data[12] | | |

ether_output

Check RT

Resolve neighbor
arpresolve()

Prepare eth header

IF_ENQUEUE(&ifp->if_snd)

end

arpresolve

Check bcast and mcast

arplookup()

Check ARP

**la_hold** update

arpwhohas()

end

---

ether_input

LINKUP?

Set MCAST/BCAST
Count the packet

Switch ethtype

ETHERTYPE_ARP → arp_intrq

**ETHERTYPE_IP** → ip_intrq

IF_ENQUEUE

end

arpintr

arpintrq.ifq_head?

IF_DEQUEUE

Is a ARP?

in_arpinput

end

---

in_arpinput

Get mac and IP from MBuf

Traverse in_ifaddr

does the src IP or dst IP Exists?

Is the packet from myself?
Src mac == outif mac

Is source MAC Bcast?

Src IP == local IP?

arp_lookup
Search for an RT entry for the src IP;
Alloc an ARP entry for src IP if needed.

Dst MAC changed?

Update DMAC using the sender MAC address

Hold packet?

Send the packet

ARP request and I am the target?

Arp proxy process

Compose arp reply,
Reply it.

end

---

ipintr

IF_DEQUEUE

Pull up IP

Is IPV4?

ip_forward()

Check DIP and BCAST,
Deliver to myself.

ip_mforward()

Reassemble packet

Demux packet
(pr_input)

udp_input    tcp_input    icmp_input    igmp_input    rip_input

---



硬件类型，ar_hrd(ARPHRD_ETHER)
协议类型，ar_pro(ETHERTYPE_IP)
硬件地址长度，ar_hln(6)
协议地址长度,ar_pln(4)

| ether_dhost | ether_shost | ether_type | | | ar_op | arp_sha | arp_spa | arp_tha | arp_tpa |
|---|---|---|---|---|---|---|---|---|---|
| 以太网目的地址 | 以太网源地址 | 帧类型 | | | op | 发送者硬件地址 | 发送者IP地址 | 目标硬件地址 | 目标IP地址 |
| 6 bytes | 6 | 2 | 2  2  1  1 | 2 | | 6 | 4 | 6 | 4 |

以太网首部    ARP首部
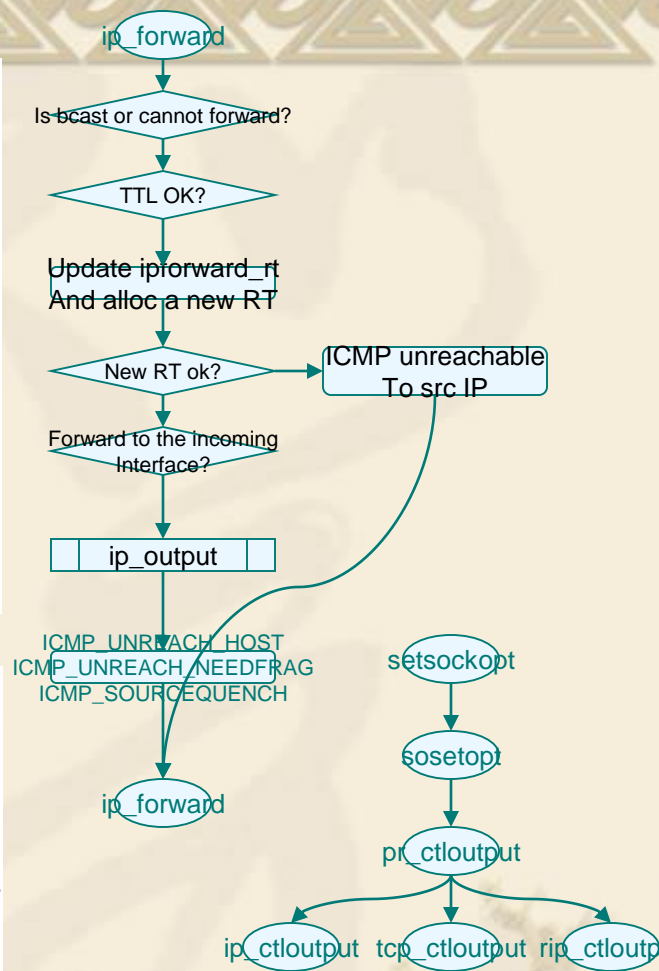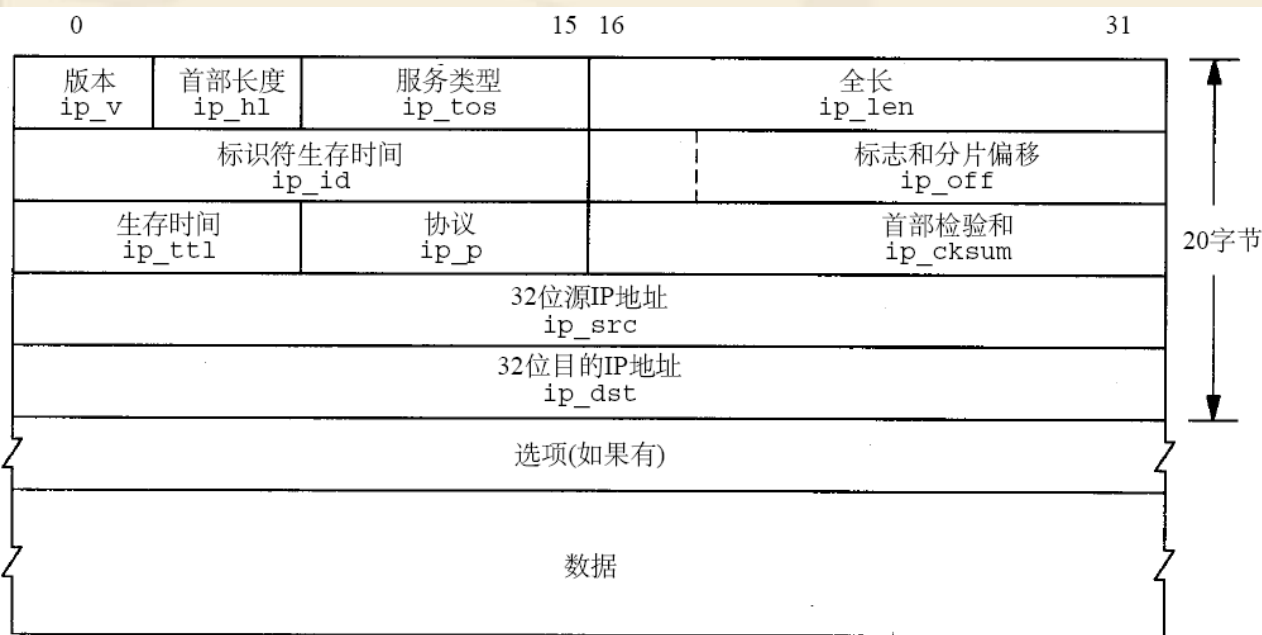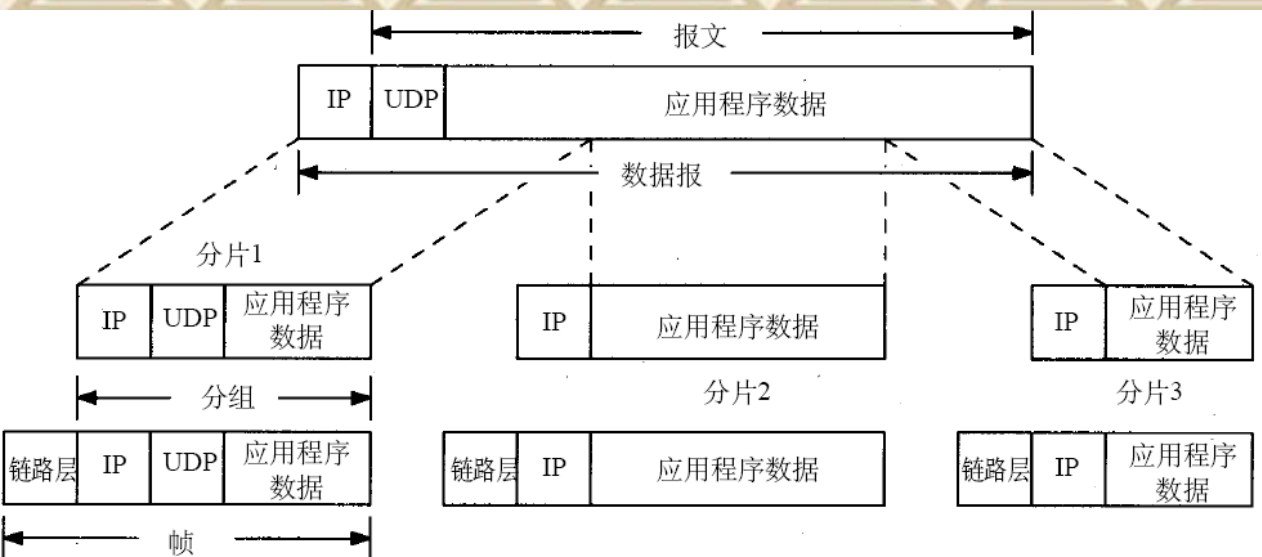ether_header()    arphdr{}    以太网ARP字段
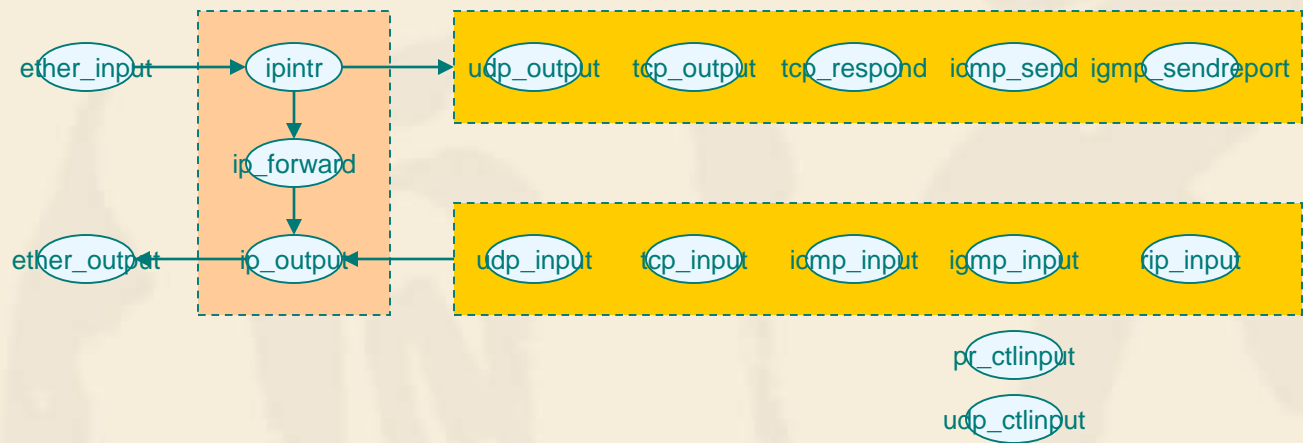ether_arp{}

图21-7 在以太网上使用时ARP请求或回答的格式

---

**TCP分段与IP分片区别:**
**1.IP分片产生的原因是网络层的MTU；TCP分段产生原因是MSS.**
**2.IP分片由网络层完成，也在网络层进行重组；TCP分段是在传输层完成，并在传输层进行重组. //透明性**
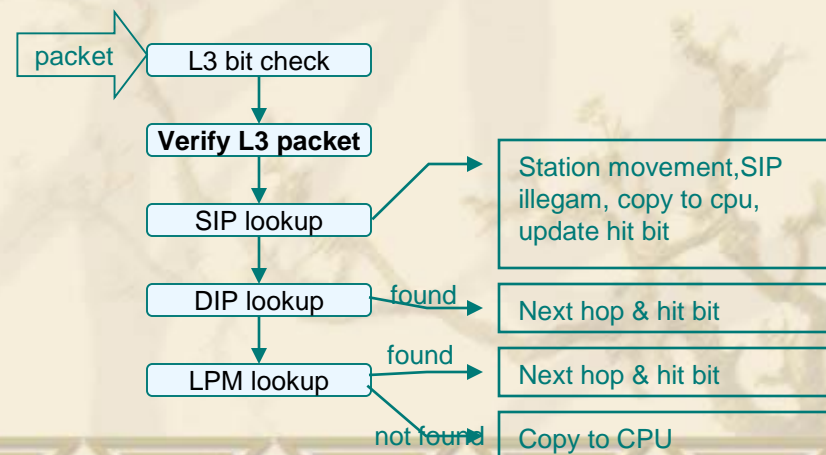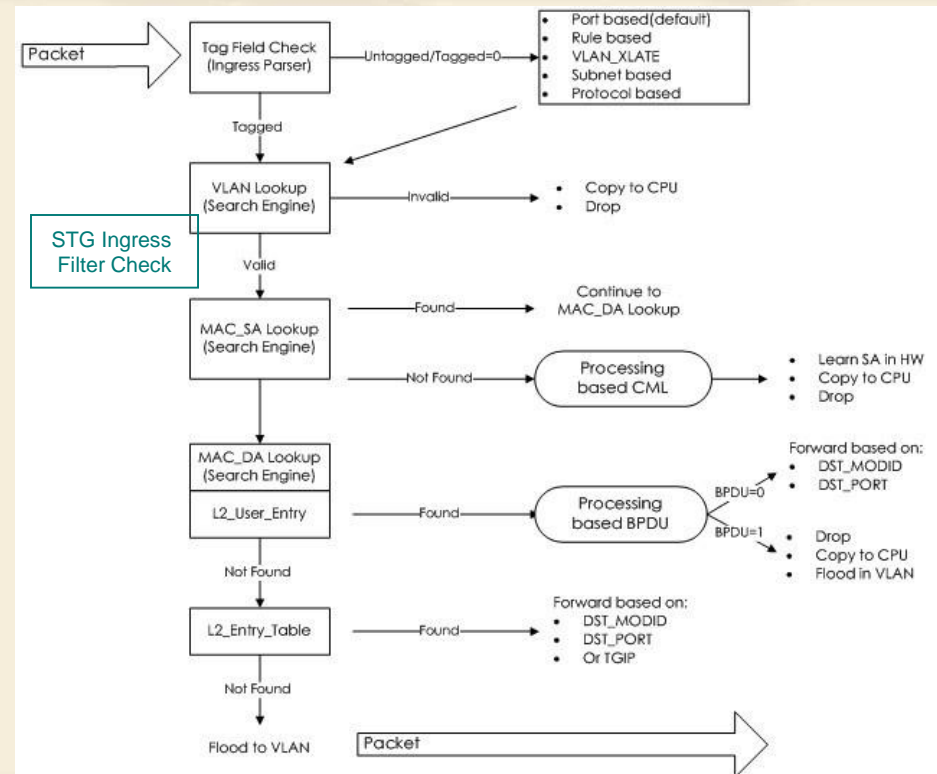**3.对于以太网，MSS为1460字节，而MTU往往会大于MSS.**

**故采用TCP协议进行数据传输，是不会造成IP分片的。若数据过大，只会在传输层进行数据分段，到了IP层就不用分片。**
**而我们常提到的IP分片是由于UDP传输协议造成的，因为UDP传输协议并未限定传输数据报的大小。**

报文

| IP | UDP | 应用程序数据 |

数据报

分片1

| IP | UDP | 应用程序数据 |

分组

| 链路层 | IP | UDP | 应用程序数据 |

帧

分片2

| IP | 应用程序数据 |

| 链路层 | IP | 应用程序数据 |

分片3

| IP | 应用程序数据 |

| 链路层 | IP | 应用程序数据 |

| 0 | 15 16 | 31 |
|---|---|---|
| 版本 ip_v | 首部长度 ip_hl | 服务类型 ip_tos | 全长 ip_len |

| 标识符生存时间 ip_id | 标志和分片偏移 ip_off |
| 生存时间 ip_ttl | 协议 ip_p | 首部检验和 ip_cksum |
| 32位源IP地址 ip_src |
| 32位目的IP地址 ip_dst |
| 选项(如果有) |
| 数据 |

20字节

**ip_forward**

Is bcast or cannot forward?

TTL OK?

Update ipforward_rt
And alloc a new RT

New RT ok? → ICMP unreachable To src IP

Forward to the Incoming Interface?

ip_output

ICMP_UNREACH_HOST
ICMP_UNREACH_NEEDFRAG
ICMP_SOURCEQUENCH

ip_forward

setsockopt

sosetopt

pr_ctloutput

ip_ctloutput   tcp_ctloutput   rip_ctloutp
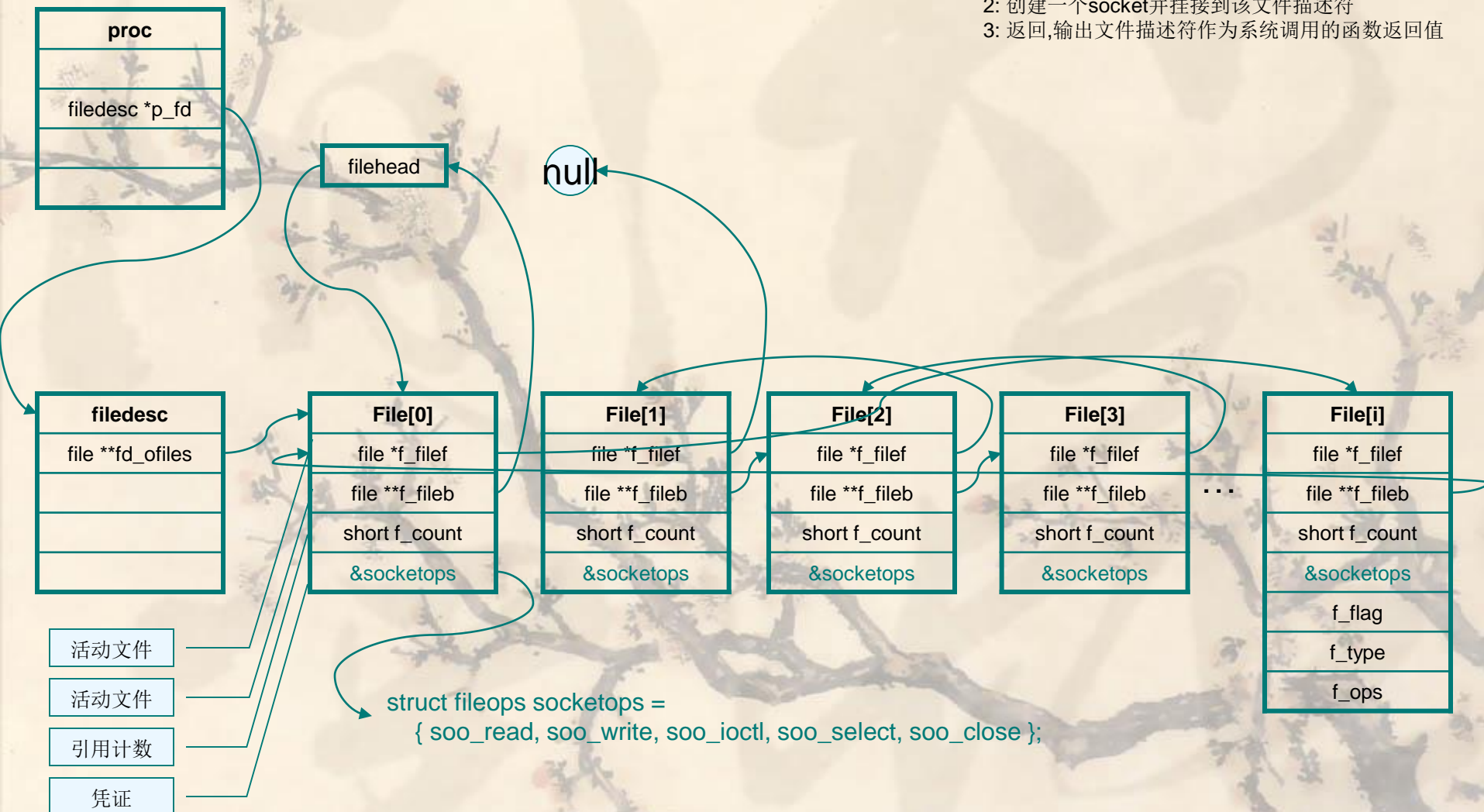
SYSCTL命令:
0, IPCTL_FORWARDING 转发功能使能
1, IPCTL_SENDREDIRECTS ICMP重定向使能
2, IPCTL_DEFTTL 设置默认TTL
3, IPCTL_DEFMTU 设置IP的MTU

STG Ingress Filter Check

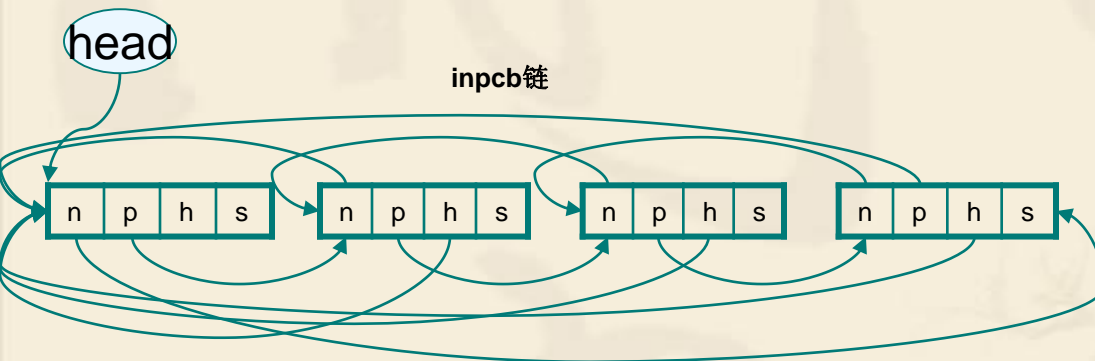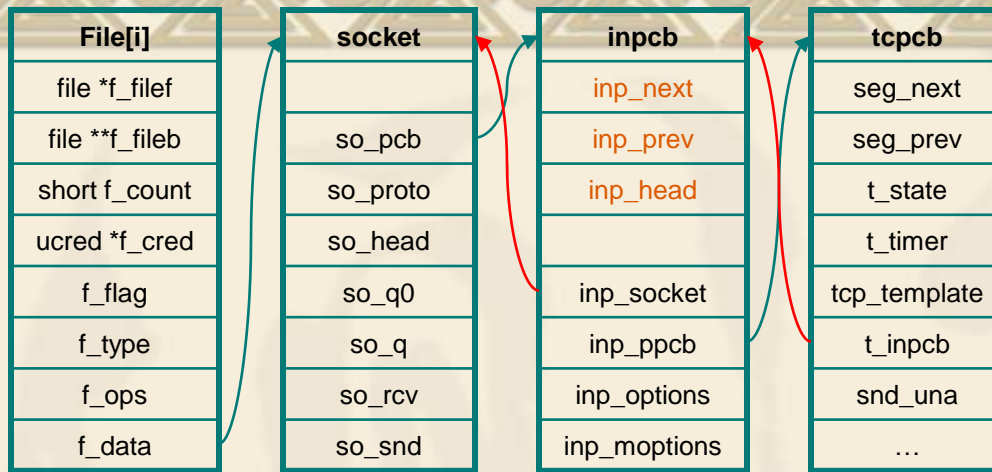

packet → L3 bit check

**Verify L3 packet**

SIP lookup → Station movement,SIP illegam, copy to cpu, update hit bit

DIP lookup — found → Next hop & hit bit

LPM lookup — found → Next hop & hit bit

not found → Copy to CPU

# socket(p,uap,retval)

1: 创建一个文件描述符并挂接到进程的活动文件列表
2: 创建一个socket并挂接到该文件描述符
3: 返回,输出文件描述符作为系统调用的函数返回值

**proc**

| |
|---|
| filedesc *p_fd |
| |
| |

filehead

null

**filedesc**

| |
|---|
| file **fd_ofiles |
| |
| |
| |

活动文件

活动文件

引用计数

凭证

| **File[0]** |
|---|
| file *f_filef |
| file **f_fileb |
| short f_count |
| &socketops |

| **File[1]** |
|---|
| file *f_filef |
| file **f_fileb |
| short f_count |
| &socketops |

| **File[2]** |
|---|
| file *f_filef |
| file **f_fileb |
| short f_count |
| &socketops |

| **File[3]** |
|---|
| file *f_filef |
| file **f_fileb |
| short f_count |
| &socketops |

· · ·

| **File[i]** |
|---|
| file *f_filef |
| file **f_fileb |
| short f_count |
| &socketops |
| f_flag |
| f_type |
| f_ops |

struct fileops socketops =
    { soo_read, soo_write, soo_ioctl, soo_select, soo_close };

falloc(p, &fp, &fd)

| File[i] |
| --- |
| file *f_filef |
| file **f_fileb |
| short f_count |
| ucred *f_cred |
| f_flag |
| f_type |
| f_ops |
| f_data |

| socket |
| --- |
| |
| so_pcb |
| so_proto |
| so_head |
| so_q0 |
| so_q |
| so_rcv |
| so_snd |

| inpcb |
| --- |
| inp_next |
| inp_prev |
| inp_head |
| |
| inp_socket |
| inp_ppcb |
| inp_options |
| inp_moptions |

| tcpcb |
| --- |
| seg_next |
| seg_prev |
| t_state |
| t_timer |
| tcp_template |
| t_inpcb |
| snd_una |
| … |

head

**inpcb链**

| n | p | h | s | | n | p | h | s | | n | p | h | s | | n | p | h | s |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**参数:**
domain: 所使用协议族
| PF_INET | TCP/IP协议; |
| PF_OSI,PF_ISO | OSI |
| PF_LOCAL,PF_UNIX | 本地IPC |
| PF_ROUTE | 路由表 |
| n/a | 链路层 |

type: socket的类型(协议的通信语义)
| SOCK_STREAM | (TCP) |
| SOCK_DGRAM(UDP) | |
| SOCK_RAW | (ICMP, IGMP, raw IP) |

protocol:
| 0 | IP协议 |
| IPPROTO_UDPUDP | |
| IPPROTO_TCPTCP | |
| IPPROTO_RAW | raw IP |
| IPPROTO_ICMP | ICMP |
| IPPROTO_IGMP | IGMP |

**过程:**
1, 基于domain, type以及protocol查找protosw;
2, 分配并初始化新的socket;重点是挂接protosw到socket

**3, 使用具体协议的usrreq进行PRU_ATTACH请求\***
3.1: in_pcballoc进行inpcb的创建和插入
3.2: soreserve进行mbuf的大小支配(此时并未实际分配内存**)**

4, 返回,输出参数为新创建的socket

socreate(uap-domain, &so, uap->type, uap->protocol)

| mbuf |
| --- |
| m_next |
| m_nextpkt |
| 0~108 |
| m_data |
| MT_xxx |
| 0 |
| 108字节dat |

| mbuf |
| --- |
| m_next |
| m_nextpkt |
| 0~100 |
| m_data |
| MT_xxx |
| M_PKTHDR |
| |
| 100字节data |

| mbuf |
| --- |
| m_next |
| m_nextpkt |
| 208~2048 |
| m_data |
| MT_xxx |
| M_EXT |
| |
| |
| |
| 2048 |

| mbuf |
| --- |
| m_next |
| m_nextpkt |
| 208~2048 |
| m_data |
| MT_xxx |
| M_PKTHDR|M_EXT |
| |
| |
| |
| 2048 |

m_next
m_nextpkt
m_len       } m_hdr{}
m_data
m_type
m_flags

m_pkthdr.len } pkthdr{}
m_pkthdr.rcvif

m_ext.ext_buf
m_ext.ext_free
m_ext.ext_size

四种不同的mbuf, 若处理小于208的
数据,可以用前两种mbuf, 若超过208
需要用到后面两种

| cluster(2048) |
| --- |
| |
| |

| cluster(2048) |
| --- |
| |
| |

listen(p,uap,retval)

**参数:**
s: socket的文件描述符
name: 传输地址或者主机名 (sockaddr)
namelen:
**过程:**
1, getsock()获取文件file指针;
2, sockargs()将传入参数复制到内核中的一个
新分配的mbuf中 (m_flags=0); mbuf存储的
sockaddr格式的数据.
**3, sobind进行usrreq请求PRU_BIND.**
3.1 in_pcbbind(inp, nam):
说明: 为inpcb设定地址和端口;
A, 接口必须有地址; 插口还未绑定;
B, socket选项是否允许本地地址复用;
C, 是否为请求连接或者接受连接
D, 是否允许通配. 匹配若B,C都真则.
wild=INPLOOKUP_WILDCARD;
E, 若为组播地址??
F, 若为单播, 查找一个地址相同的接口
G, 端口不可为保留端口
H, 根据端口和地址查找可用pcb,
I, 判断该pcb是否允许复用端口和地址
J, 若未指明绑定端口号,则从pcb链中查找一个
可用的端口号分配给用户使用.

**参数:**
s: socket的文件描述符,指明哪个socket被设为被
动模式;
backlog: 指明socket使用队列的长度
**过程:**
1, getsock()获取文件file指针;
**2, 进行usrreq请求PRU_LISTEN.**
2.1 in_pcbbind(inp, nam):
说明: 见bind
UDP不支持listen
与udp不同处:
A, inpcb转化成tcpcb
B, tcpcb设置为监听状态.
3, 设置sock模式为SO_ACCEPTCONN, 设置
sock的队列最大长度为backlog.

# accept(p,uap,retval)
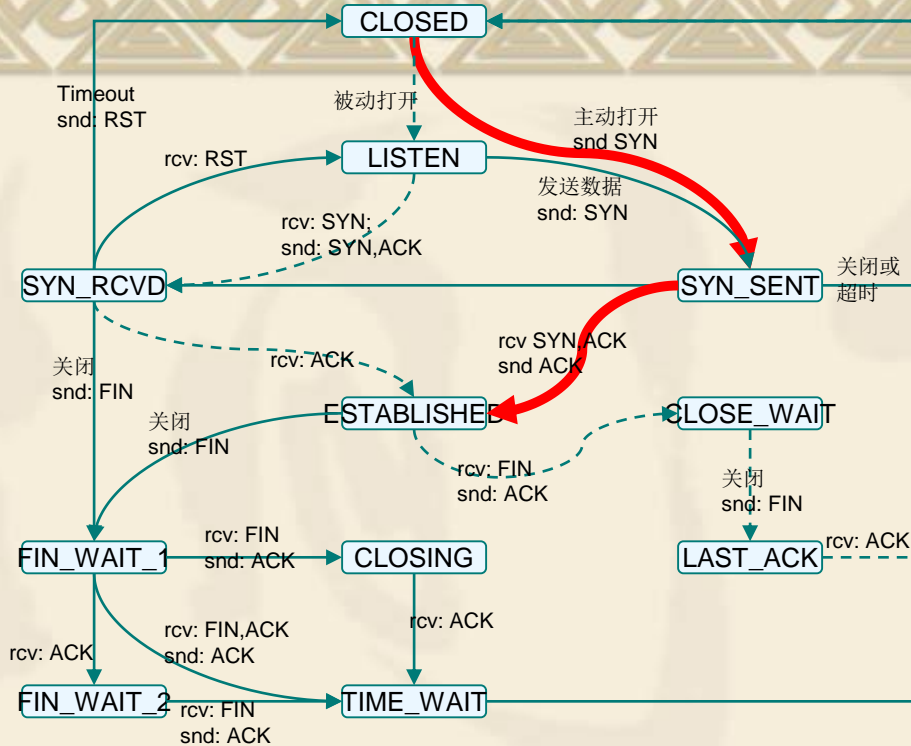
## TCP状态图

CLOSED

Timeout
snd: RST

被动打开

LISTEN

rcv: SYN;
snd: SYN,ACK

SYN_RCVD          SYN_SENT

rcv: ACK

ESTABLISHED       CLOSE_WAIT

rcv: FIN
snd: ACK

关闭
snd: FIN

FIN_WAIT_1        CLOSING           LAST_ACK    rcv: ACK

FIN_WAIT_2        TIME_WAIT

## 流程图

校验,保存参数

取得sock

SO_ACCEPTCONN?
↓ N

SO_NBIO?
↓ N

so_qlen==0&&so_error==0?
↓ Y

SS_CANTRCVMORE?
↓ N

挂起等待连接

falloc一个file

soqremove
从接收队列中删除
(TCP状态机插入)

将sock放入新的file;
新分配一个mbuf:
nam

PRE_ACCEPT请求
nam存储对端传输地址

复制nam和namelen
给进程,释放nam

## 队列结构

so_head
so_q0
so_q

so_head / so_q0 / so_q    so_head / so_q0 / so_q    null

so_head / so_q0 / so_q    so_head / so_q0 / so_q    so_head / so_q0 / so_q

## 说明

**参数:**
s: socket的文件描述符,指明哪个socket被设为被动模式;
name: 缓冲指针,用来保存外部主机传输地址
anamelen: 缓冲指针长度指针
**过程:**
1, getsock()获取文件file指针;
2, **进行usrreq请求PRU_LISTEN.**
2.1 in_pcbbind(inp, nam):
说明: 见bind
UDP不支持listen,accept
与udp不同处:
A, inpcb转化成tcpcb
B, tcpcb设置为监听状态.
3, 设置sock模式为SO_ACCEPTCONN, 设置sock的队列最大长度为backlog.

状态变迁说明:
1. 当服务器状态由LISTEN到SYN_RCVD变迁的时候, 新进来的连接将会使tcp_input产生一个新的sock,并将该sock插入到半连接队列.

2. 只有当服务器状态由SYN_RCVD变迁到ESTABLISHED的时候进入tcp_input的SYN_RCVD分支, soisconnected()函数被调用, 挂起的accept才被唤醒并将sock插入到连接队列

PRE_ACCEPT:
直接通过mbuf返回对端传输地址

# 状态转换图

CLOSED

Timeout
snd: RST

被动打开

主动打开
snd SYN

rcv: RST

LISTEN

发送数据
snd: SYN

rcv: SYN;
snd: SYN,ACK

SYN_RCVD

SYN_SENT

关闭或
超时

rcv SYN,ACK
snd ACK

关闭
snd: FIN

rcv: ACK

ESTABLISHED

CLOSE_WAIT

关闭
snd: FIN

rcv: FIN
snd: ACK

关闭
snd: FIN

FIN_WAIT_1

rcv: FIN
snd: ACK

CLOSING

LAST_ACK

rcv: ACK

rcv: ACK

rcv: FIN,ACK
snd: ACK

rcv: ACK

rcv: ACK

FIN_WAIT_2

rcv: FIN
snd: ACK

TIME_WAIT

## connect(p,uap,retval)

**参数:**
s: socket的文件描述符;
name: 指向远端地址
namelen: name的长度
connect既可用在TCP也可用在UDP;用于主动初始化一条连接,为客户端函数.
**过程:**
1, getsock()获取文件file指针;
2, sock状态不能为NBIO和CONNECTING;
3, sockargs()将传入参数复制到内核中的一个新分配的mbuf中 (m_flags=0);
**4, soconnect进行usrreq请求PRU_LISTEN.**
UDP: in_pcbconnect
A, 校验和转换外部地址
B, 路由判断,首次发送需要建立一条路由并删除老的路由,再次发送可以直接使用已有路由.
C, 选择路由出接口对应的本地地址()
D, in_pcblookup验证插口对的唯一性
E, in_pcbbind隐式绑定本地端口
F, 分配远端地址和端口

TCP:
1, 若为bind则调用in_pcbbind进行本地端口绑定
2, 同UDP, 使用in_pcbconnect对sock指定本地IP和外部端口号
3, tcp_template为连接的IP和TCP首部创建一个模板
**4, 计算窗口缩放因子**
5, sock设为connecting状态, tcp_connattempt++
6, TCP状态机进入SYN_SENT状态,
7, 初始化TCP报文序号,
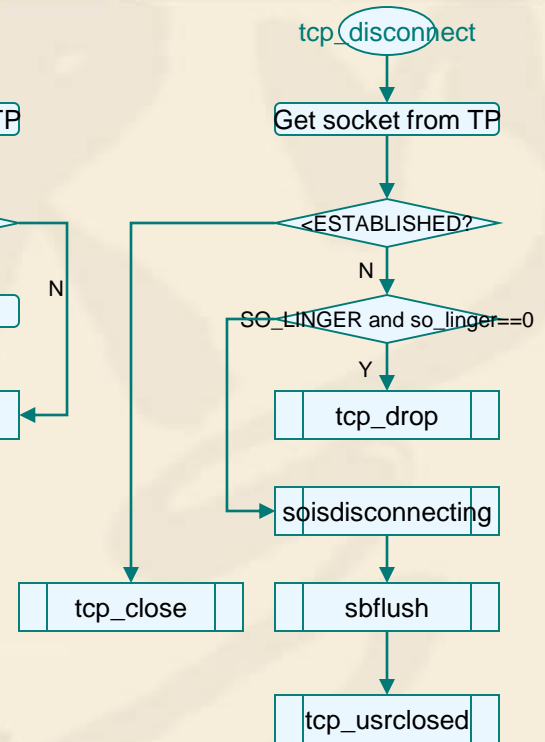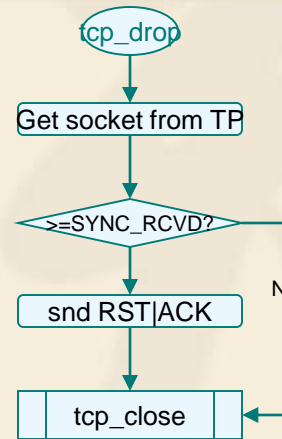8, tcp_output发送TCP报文,此时由于处于SYN_SENT, 发送SYN报文到服务器
9, tsleep等待连接建立, 并取消sock的connecting状态.
10, 从服务器端收到对应的ACK, 发送ACK并进入ESTABLISHED状态并抹掉sock的CONNECTING状态.(此时CONNECTED被TCP状态机置位)

## shutdown(p,uap,retval)

**参数:**
s: socket的文件描述符;
how: 连接关闭的方式
FREAD关闭连接的读通道,FWRITE关闭连接的写通道,FREAD|FWRITE关闭连接的读写通道.
**过程:**
1, getsock()获取文件file指针;
2, usrreq-SHUTDOWN
FREAD: sock设置为CANTRCVMORE唤醒进程;
对于注册了dispose的协议,调用dispose释放资源
FWRITE: sock设置为CANTSENDMORE唤醒进程;

CLOSED

LISTEN

SYN_RCVD

SYN_SENT

ESTABLISHED

CLOSE_WAIT

FIN_WAIT_1

CLOSING

LAST_ACK

FIN_WAIT_2

TIME_WAIT

Timeout
snd: RST

被动打开

rcv: RST

主动打开
snd SYN

发送数据
snd: SYN

rcv: SYN;
snd: SYN,ACK

关闭或
超时

rcv: ACK

rcv SYN,ACK
snd ACK

关闭
snd: FIN

关闭
snd: FIN

rcv: FIN
snd: ACK

关闭
snd: FIN

rcv: FIN
snd: ACK

rcv: ACK

rcv: FIN,ACK
snd: ACK

rcv: ACK

rcv: FIN
snd: ACK

rcv: ACK

## tcp_drop

Get socket from TP

\>=SYNC_RCVD?

snd RST|ACK

tcp_close

N

tcp_close

## tcp_disconnect

Get socket from TP

<ESTABLISHED?

N

SO_LINGER and so_linger==0

Y

tcp_drop

soisdisconnecting

sbflush

tcp_usrclosed

## tcp_close

Update some info in
inp_route if necessary

free the reass queue

Free the tcpcb

Remove the sock from queue,
And free it.

Remove the inpcb from queue,
And free it.

# close(p,uap,retval)

**参数:**
s: socket的文件描述符;
**过程:**
0, 对于未连接的sock, usrreq-ABORT 释放
q,q0
1, 对于已连接的sock, usrreq-DISCONNECT
2, 如果是connected状态则挂起等待
3, usrreq-DETACH,
4, 释放sock
5, 释放file

| File[i] |
|---|
| file *f_filef |
| file **f_fileb |
| short f_count |
| ucred *f_cred |
| f_flag |
| f_type |
| f_ops |
| f_data |

| socket |
|---|
| |
| so_pcb |
| so_proto |
| so_head |
| so_q0 |
| so_q |
| so_rcv |
| so_snd |

| inpcb |
|---|
| inp_next |
| inp_prev |
| inp_head |
| |
| inp_socket |
| inp_ppcb |
| inp_options |
| inp_moptions |

| tcpcb |
|---|
| seg_next |
| seg_prev |
| t_state |
| t_timer |
| tcp_template |
| t_inpcb |
| snd_una |
| … |

| sockbuf |
|---|
| sb_mb |
| //snd queue |

| mbuf |
|---|
| m_next |
| m_nextpkt |
| |

read
Read data to recv queue;
recv
Read data to recv queue, with options
write
Write data to snd queue
send
Write data to snd queue, with options
select
Waiting for I/O events

#define TH_FIN    0x01 //finalize
#define TH_SYN    0x02 //sync, init a connection
#define TH_RST    0x04 //reset
#define TH_PUSH  0x08 //push to app
#define TH_ACK    0x10 //ACK on
#define TH_URG    0x20 //urgent pointer on

#define TCP_MSS 512 //Maximum Segment Size, to avoid fragmentation

#define TCPTV_MSL ( 30*PR_SLOWHZ)
//max seg lifetime.
Usage: 1, When TCP is in TIME_WAIT state, if receives a RST, drop the connection. If it receives another FIN, the timer should be reset and wait for another 2MSL.
2, Quiet time. When TCP reset, no connection can be established in MSL.

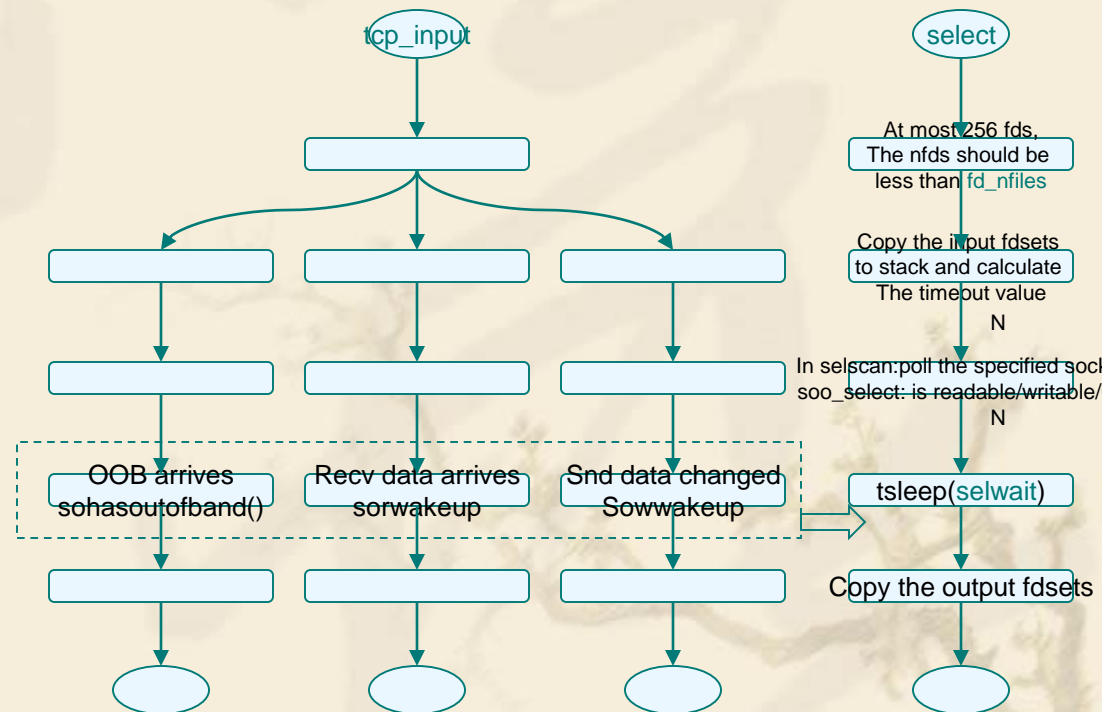**select(int nfds, fd_set \*readfds, fd_set \*writefds, fd_set \*exceptfds, struct timeval \*timeout)**
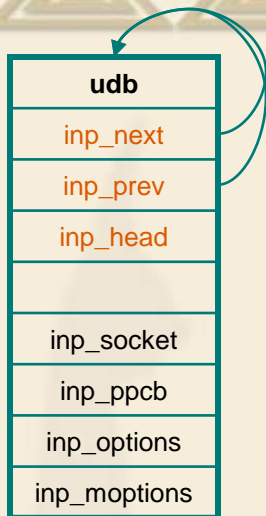
**Params:**
nfds: max handle
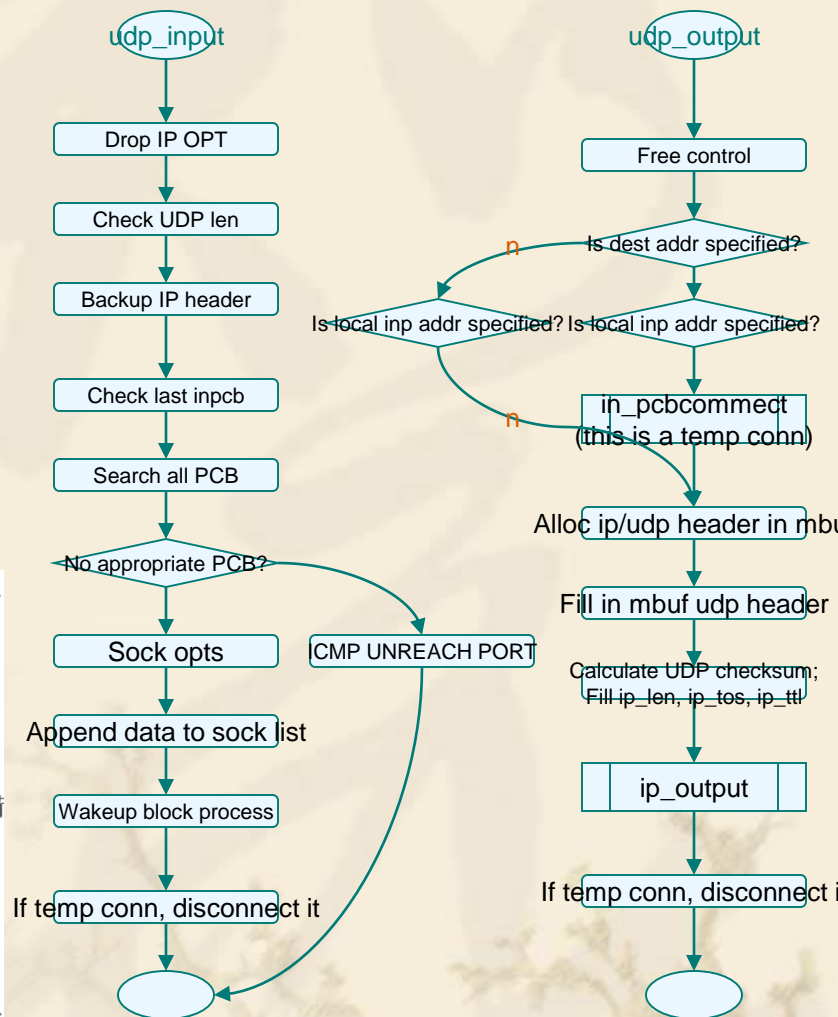fd_set: three types of monitored fd;
timeout:

Compare with poll/epoll

tcp_input

select

At most 256 fds,
The nfds should be
less than fd_nfiles

Copy the input fdsets
to stack and calculate
The timeout value
N

In selscan:poll the specified sock
soo_select: is readable/writable/(
N

OOB arrives
sohasoutofband()

Recv data arrives
sorwakeup

Snd data changed
Sowwakeup

tsleep(selwait)

Copy the output fdsets

**udb**

| inp_next |
| inp_prev |
| inp_head |
| |
| inp_socket |
| inp_ppcb |
| inp_options |
| inp_moptions |

**udp_init()**

SYSCTL命令:
0, IPCTL_FORWARDING 转发功能使能
1, IPCTL_SENDREDIRECTS ICMP重定向使能
2, IPCTL_DEFTTL 设置默认TTL
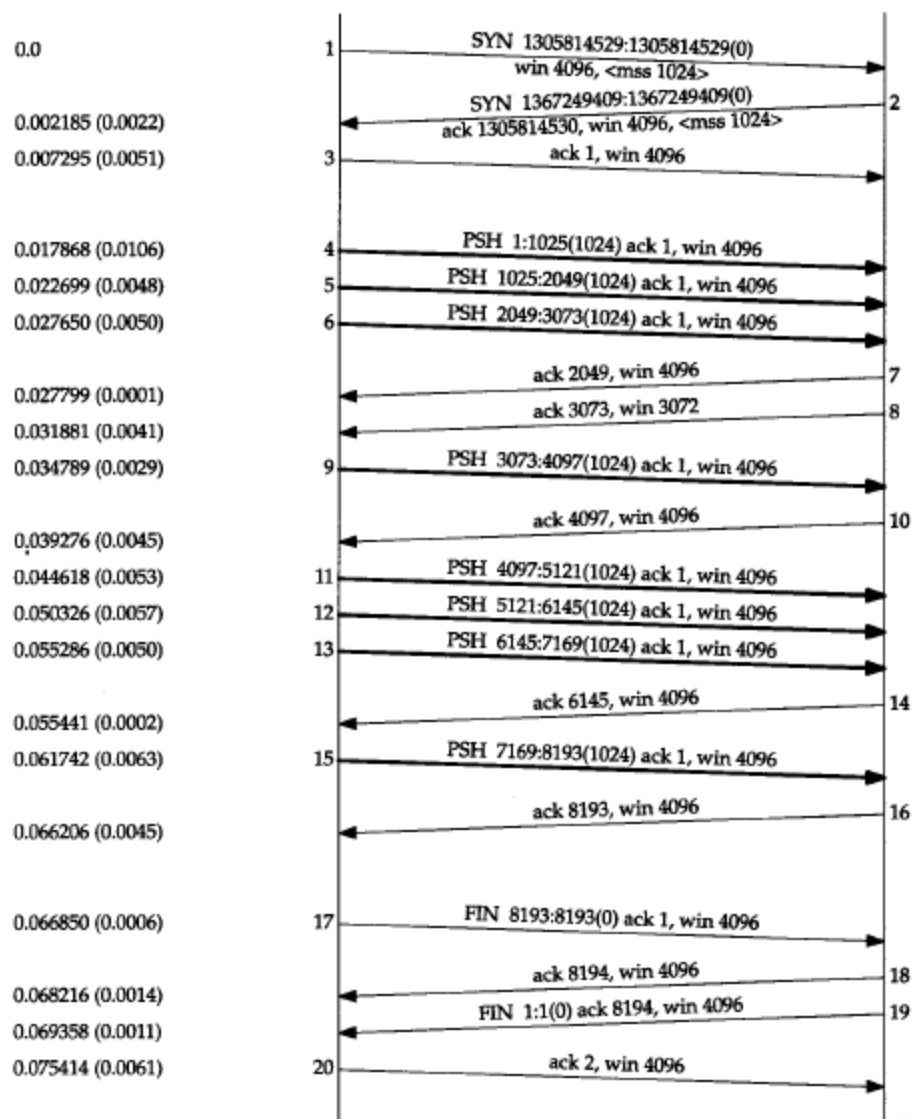3, IPCTL_DEFMTU 设置IP的MTU

4, UDPCTL_CHECKSUM 使能UDP检验和

| 0 | 15 | 16 | 31 |
|---|---|---|---|
| 4位版本号 4位首部长度 | 8位服务类型 (TOS) | 16位全长(字节数) | |
| 16位标识符 | 3位标志 | 13位分片偏移 | |
| 8位生存时间 (TTL) | 8位协议 | 16位首部检验和 | |
| 32位源IP地址 | | | |
| 32位目的IP地址 | | | |
| 16位源端口号 | | 16位目的端口号 | |
| 16位UDP长度 | | 16位UDP检验和 | |

20字节

8字节

**udp_input**

- Drop IP OPT
- Check UDP len
- Backup IP header
- Check last inpcb
- Search all PCB
- No appropriate PCB?
- Sock opts
- Append data to sock list
- Wakeup block process
- If temp conn, disconnect it

CMP UNREACH PORT

**udp_output**

- Free control
- Is dest addr specified?
- Is local inp addr specified? Is local inp addr specified?
- in_pcbcommect (this is a temp conn)
- Alloc ip/udp header in mbuf
- Fill in mbuf udp header
- Calculate UDP checksum; Fill ip_len, ip_tos, ip_ttl
- ip_output
- If temp conn, disconnect it

n

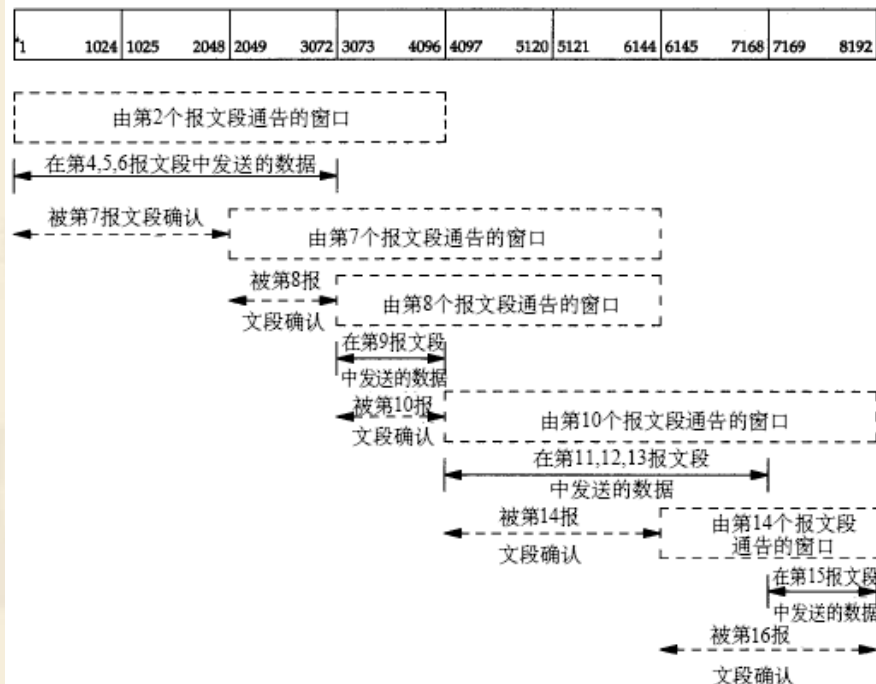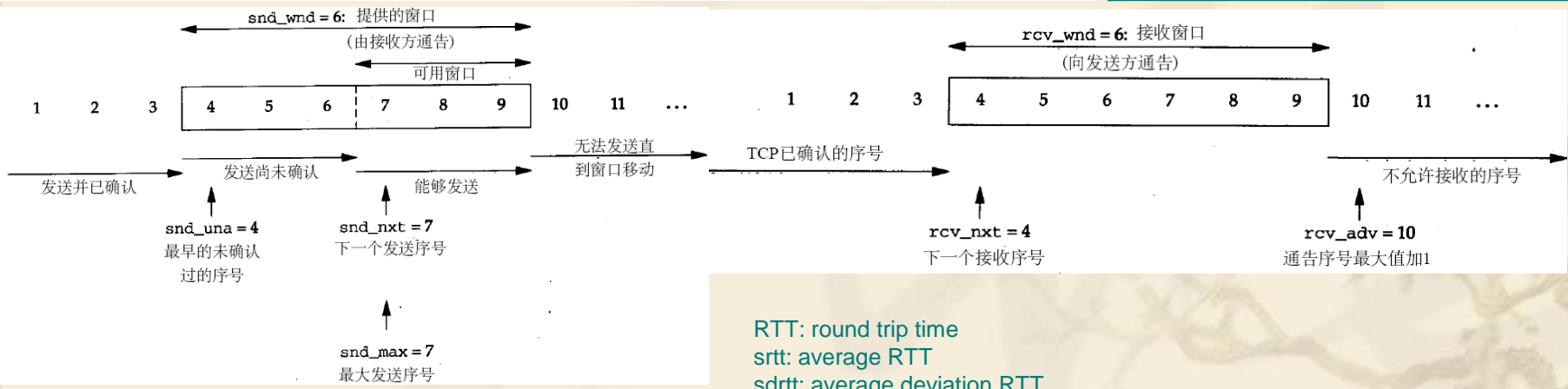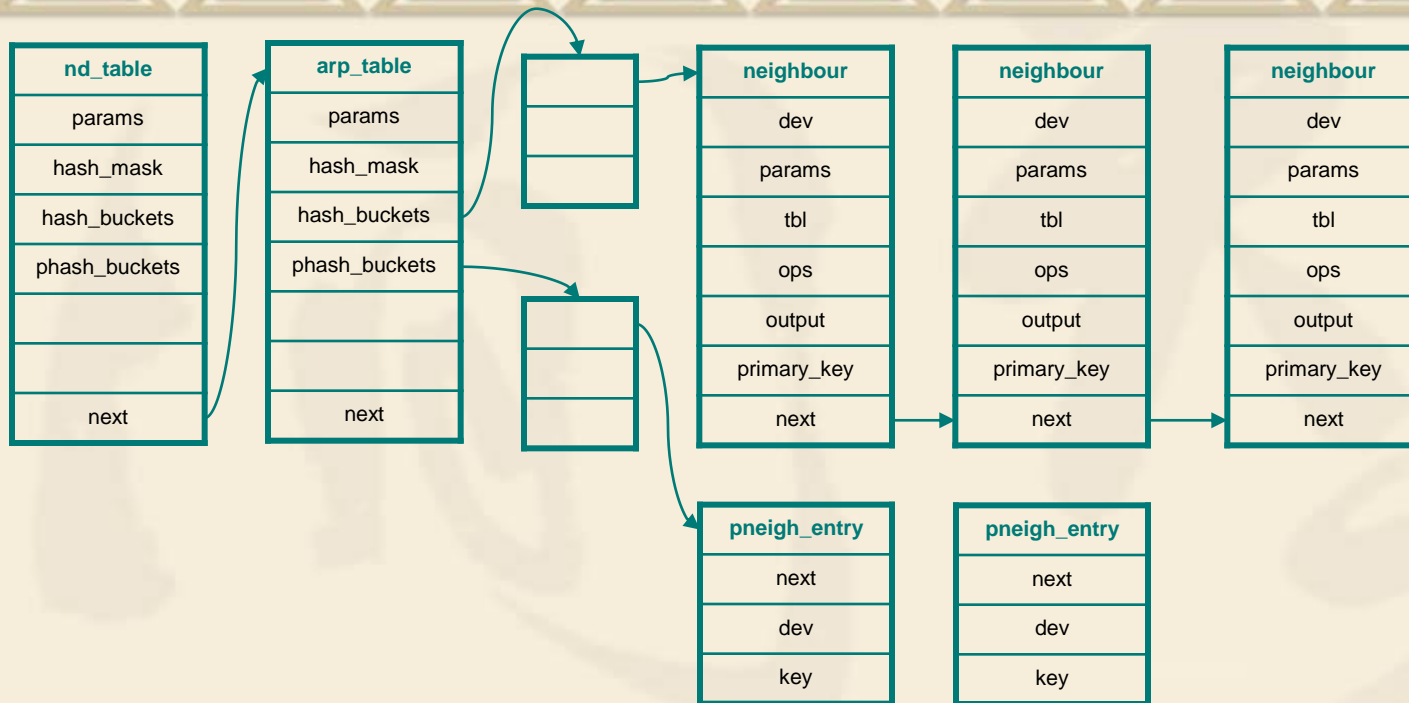| name | Description |
|------|-------------|
| snd_una | The first unacknowledged no. |
| snd_nxt | The next transmitted no. |
| snd_wnd | The offered window size |
| snd_max | The maximum sent no. (snd_nxt<snd_max, retransmitting) |
| rcv_wnd | Recv window to inform the sender |
| rcv_nxt | The next recv no. |
| rcv_adv | the first no. outside window |
| | |
| | |

svr4.1056                                                      bsdi.7777

0.0                    1    SYN  1305814529:1305814529(0)
                            win 4096, <mss 1024>
                            SYN  1367249409:1367249409(0)      2
0.002185 (0.0022)           ack 1305814530, win 4096, <mss 1024>
0.007295 (0.0051)      3    ack 1, win 4096

0.017868 (0.0106)      4    PSH  1:1025(1024) ack 1, win 4096
0.022699 (0.0048)      5    PSH  1025:2049(1024) ack 1, win 4096
0.027650 (0.0050)      6    PSH  2049:3073(1024) ack 1, win 4096

0.027799 (0.0001)           ack 2049, win 4096                 7
0.031881 (0.0041)           ack 3073, win 3072                 8
0.034789 (0.0029)      9    PSH  3073:4097(1024) ack 1, win 4096

0.039276 (0.0045)           ack 4097, win 4096                 10
0.044618 (0.0053)      11   PSH  4097:5121(1024) ack 1, win 4096
0.050326 (0.0057)      12   PSH  5121:6145(1024) ack 1, win 4096
0.055286 (0.0050)      13   PSH  6145:7169(1024) ack 1, win 4096

0.055441 (0.0002)           ack 6145, win 4096                 14
0.061742 (0.0063)      15   PSH  7169:8193(1024) ack 1, win 4096

0.066206 (0.0045)           ack 8193, win 4096                 16

0.066850 (0.0006)      17   FIN  8193:8193(0) ack 1, win 4096

0.068216 (0.0014)           ack 8194, win 4096                 18
0.069358 (0.0011)           FIN  1:1(0) ack 8194, win 4096     19
0.075414 (0.0061)      20   ack 2, win 4096

图20-1  从svr4 传输8192个字节到bsdi

| 1 | 1024 | 1025 | 2048 | 2049 | 3072 | 3073 | 4096 | 4097 | 5120 | 5121 | 6144 | 6145 | 7168 | 7169 | 8192 |

由第2个报文段通告的窗口

在第4,5,6报文段中发送的数据

被第7报文段确认

由第7个报文段通告的窗口

被第8报
文段确认

由第8个报文段通告的窗口

在第9报文段
中发送的数据

被第10报
文段确认

由第10个报文段通告的窗口

在第11,12,13报文段
中发送的数据

被第14报
文段确认

由第14个报文段
通告的窗口

在第15报文段
中发送的数据

被第16报
文段确认

图20-6  图20-1的滑动窗口协议

## TCP Header

| th_sport 16位源端口号 | th_dport 16位目的端口号 |
|---|---|
| th_seq 32位序号 | |
| th_ack 32位确认序号 | |

| th_off 4位 首部长度 | th_x2 保留 (6位) | U R G | A C K | P S H | R S T | S Y N | F I N | th_win 16位窗口大小 |

| th_sum 16位TCP检验和 | th_urp 16位紧急数据偏移量 |

选项(如果有)

数据(如果有)

20字节

| timer | description |
|---|---|
| connection | SYN, 75s |
| **retransmission** | **Dynamic, retransmit data** |
| Delayed ACK | 捎带确认 ACK doesn't need to be confirmed immediately. (fast timer) |
| **persist** | **Win=0, send 1 byte, dynamic** **Test the availability of the peer's window size.** |
| keepalive | SO_KEEPALIVE, alive test |
| FIN_WAIT_2 | 10m+75s, close called Avoid the losing FIN from peer. |
| TIME_WAIT | 2 MSL |

snd_wnd = 6: 提供的窗口 (由接收方通告)
可用窗口

1 2 3 | 4 5 6 | 7 8 9 | 10 11 ...

发送并已确认
发送尚未确认
能够发送
无法发送直 到窗口移动

snd_una = 4 最早的未确认 过的序号
snd_nxt = 7 下一个发送序号
snd_max = 7 最大发送序号

rcv_wnd = 6: 接收窗口 (向发送方通告)

1 2 3 | 4 5 6 7 8 9 | 10 11 ...

TCP已确认的序号
不允许接收的序号

rcv_nxt = 4 下一个接收序号
rcv_adv = 10 通告序号最大值加1

RTT: round trip time
srtt: average RTT
sdrtt: average deviation RTT

慢启动和拥塞避免算法
指数退避 exponential backoff

# Neighbor system(Linux)

| nd_table |
|---|
| params |
| hash_mask |
| hash_buckets |
| phash_buckets |
|  |
|  |
| next |

| arp_table |
|---|
| params |
| hash_mask |
| hash_buckets |
| phash_buckets |
|  |
|  |
| next |

| neighbour |
|---|
| dev |
| params |
| tbl |
| ops |
| output |
| primary_key |
| next |

| neighbour |
|---|
| dev |
| params |
| tbl |
| ops |
| output |
| primary_key |
| next |

| neighbour |
|---|
| dev |
| params |
| tbl |
| ops |
| output |
| primary_key |
| next |

| pneigh_entry |
|---|
| next |
| dev |
| key |

| pneigh_entry |
|---|
| next |
| dev |
| key |

struct neigh_table
struct neigh_statistics
struct neighbour
struct hhcache
struct pneigh_entry
struct neigh_parms

Memory columns (labeled 1-5), each showing addresses: 0xa0000, 0x90200, 0x90000, 0x10000, 0x07c00, 0x00000

```
0034:   SETUPLEN = 4
0035:   BOOTSEG  = 0x07c0
0036:   INITSEG  = 0x9000
0037:   SETUPSEG = 0x9020
0038:   SYSSEG   = 0x1000
0039:   ENDSEG   = SYSSEG + SYSSIZE
0040:
0041:   ! ROOT_DEV:       0x000 - same type of
0042:   !                 0x301 - first partit
0043:   ROOT_DEV = 0x306
0044:
0045:   entry start
0046:   start:
0047:           mov     ax,#BOOTSEG
0048:           mov     ds,ax
0049:           mov     ax,#INITSEG
0050:           mov     es,ax
0051:           mov     cx,#256
0052:           sub     si,si
0053:           sub     di,di
0054:           rep
0055:           movw
0056:           jmpi    go,INITSEG
0057:   go:     mov     ax,cs
0058:           mov     ds,ax
0059:           mov     es,ax
0060:   ! put stack at 0x9ff00.
0061:           mov     ss,ax
0062:           mov     sp,#0xFF00
0063:
0064:   ! load the setup-sectors directly af
0065:   ! Note that 'es' is already set up.
0066:
0067:   load_setup:
0068:           mov     dx,#0x0000
0069:           mov     cx,#0x0002
0070:           mov     bx,#0x0200
0071:           mov     ax,#0x0200+SETUPLEN
0072:           int     0x13
0073:           jnc     ok_load_setup
0074:           mov     dx,#0x0000
0075:           mov     ax,#0x0000
0076:           int     0x13
0077:           j       load_setup
0078:
```

**Diagram labels (left column boxes):**

- state
- thread_info
- run_list
- tasks
- mm
- real_parent
- parent
- tty
- thread
- fs
- files
- signal
- pending
- thread

thread_info

thread_struct

mm_struct

tty_struct

fs_struct

files_struct

signal_struct

task_struct
tasks

task_struct
tasks

task_struct
tasks

schedule

Global page directory switch

Hardware context switch

switch_to

**进程描述符task_struct**
状态、基本信息、内存描述符指针、tty、目录、文件、信号 等
**pids[PIDTYPE_MAX]: 四个散列表**
run_list
状态：TASK_RUNNING, TASK_INTERRUPTIBLE,
TASK_UNINTERRUPTIBLE, TASK_STOPPED, TASK_TRACED,
EXIT_ZOMBIE, EXIT_DEAD
设置指定进程状态 set_task_state
设置当前进程状态 set_current_state
**进程标志符** 最大为32767 设置这个最大值：/proc/sys/kernel/pid_max
Getpid & getppid
Alloc_thread_info free_thread_info;
current_thread_info: 由esp获取当前thread_info结构
**四个散列表** PIDTYPE_PID, PIDTYPE_TGID, PIDTYPE_PGID,
PIDTYPE_SID
冲突双向链表

**EFLAGS(program status and control)**程序的状态及控制

**等待队列头的创建**宏DECLARE_WAIT_QUEUE_HEAD(name)
初始化动态分配的等待队列的头变量init_waitqueue_head()
初始化wait_queue_t结构：init_waitqueue_entry(q,p)
**非互斥进程**唤醒：default_wake_function
DEFINE_WAIT声明一个wait_queue_t；autoremove_wake_function初
始化这个项；init_waitqueue_func_entry定义唤醒函数
add_wait_queue把一个非互斥进程插入到队列头；
add_wait_queue_exclusive把一个互斥进程插入到队列尾
waitqueue_active判断等待队列是否空；remove_wait_queue从等待队列
删除一个进程
**重要函数** sleep_on(), interruptible_sleep_on(), sleep_on_timeout(),
interruptible_sleep_on_timeout()
**schedule**, schedule_timeout, wait_event, wait_event_interruptible,
prepare_to_wait, prepare_to_wait_interruptible

**Kernel mode stack**: is combined with current thread info, so the
change of the stack, ie, the change of esp, means the change of
process.

fork()

ZOMBIE

schedule()

exit()

RUNNING

RUNNING

preempt

Resource available

Wait for resource

INTERRUPTIBLE
UNINTERRUPTIBLE

$(N+1)*(2<<13)-1$

stack

esp

task

task_info

task_struct

thread_info

$N*(2<<13)$

current

child.next

child.prev

p0

p1

sib.next

p2

p3

sib.prev

p4

Process
do_fork():
create a process by coping current process, share the same process context with its parent.
Return twice with different value.

do_exec(): read and load bin, alloc new context and excute it.

Thread
In linux, it is the same as process;
In Windows, it is called 'lightweight process' and supported differently with process.

**Kernel threads**: background thread, only run in kernel mode, mm is NULL.; kernel_thread().
0, 1, keventd; kapmd; kswapd; pdflush; kblockd; ksoftirqd;

**exit_group**() & **do_exit**(): exit() & pthread_exit().

Schedule
Cooperative multitasking & Preemptive multitasking
Runable task queue
Priority Array: 140

priority inversion

## Flowchart column 1: Cpu interrupt process

Cpu interrupt process

- Check int,
- Decide the int vec i,
- Read idtr[i]
- Use segment select in idt, Retrive the base addr of the entry,
- k the priviledge of the interrupt ...
- Save context in stack
- Load cs, eip from idt[i]
- Jump to handler

Cpu interrupt ret

- Restore context
- Restore ss and esp
- Check kernel preemption
- Check reschedule
- …

## Flowchart column 2: Exception call

Exception call

- Push handler,
- Push regs,
- Clr DF, copy errcode to DX,
- Load handler to DI, Write ES to the same addr
- Stack top is the addr of the hdlr, Copy it to EAX,
- Copy data segment to Ds and es,
- Call hdlr.

## Flowchart column 3: do_XXhdlr

do_XXhdlr

- current->thread.error_code current->thread.trapno
- force_sig
- …
- die(), Jmp to ret_from_exception()

## softirq_vec

```
0                    32
softirq_vec [  |  |  |  … |  ]
              |
          softirq_action
          [          ]
          [          ]
```

## Exception:
Generated by CPU itself and only occurs after an instruction finished.
1. Processor detecetd – fault, trap, abort;
2. Programmed – software interrupt, int.

## Interrupt:
Externally generated by other devices.
1. Maskable interrupt: a masked interrupt will be ignored by CPU.
2. Nonmaskable interrupt: cannot be masked.
3. 3 types of interrupt: I/O, clock, inter-CPU,
Shared IRQ: how to differentiate different IRQs coming from one IRQ line. All ISRs will be excuted, but only one will finish its function.
Dynamic IRQ: the same IRQ line visited by different devices in different times.

## Steps of I/O:
1, save IRQ and regs in stack;
2, answer to PIC;
3, ISR;
4, return from intr.

## Inter-Processor Interrupt: (call function, reschedule, invalidata TLB)
1, save IRQ and regs in stack;
2, answer to PIC;
3, ISR;
4, return from intr.

## SoftIRQ/tasklet
**SoftIRQ:** 可重入**ksoftirqd**
**Tasklet:** 不可重入,特殊类型的软中断.

工作者线程 **worker thread**

可延迟中断
## Programmable Interrupt Controller
1,Convert IRQ to interrupt vector, store in PIC IO port;2,Send interrupt to CPU INTR;3,Wait and confirm, clear INTR.
## IDT&idtr
lidtr, 3 types of IDT defined by Intel:(task gate, interrupt gate, trap gate)
Kernel woops?
## Page fault
The only exception that can happen in kernel mode.

```
struct hw_interrupt_type i8259A_irq_type = {
    .typename     = "XT-PIC",
    .startup      = startup_8259A_irq,
    .shutdown     = shutdown_8259A_irq,
    .enable       = enable_8259A_irq,
    .disable      = disable_8259A_irq,
    .ack          = mask_and_ack_8259A,
    .end          = end_8259A_irq,
    .set_affinity = NULL
};
```
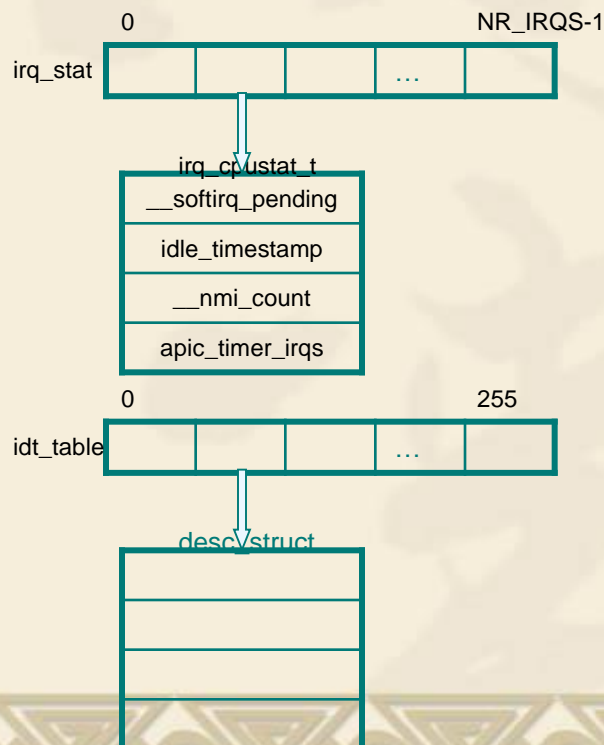
## Kernel stack:
1. Current process kernel stack used for all;
2. Exception stack, per process;
3. Hard interrupt stack. hardirq_stack[], per CPU;
4. Soft interrupt stack. softirq_stack[],per CPU;

1,所有hdlr响应,不允许产生相同的中断事件
2,中断处理,软中断和tasklet既不能被抢占也不能被堵塞
3,hdlr不能执行可延迟函数或系统调用服务的内核控制路径中断??
4,软中断和tasklet不能在一个CPU上交错执行
5,同一个tasklet不能同时在几个CPU上执行

kirqd:adjuct the CPU allocation of IRQs

Regenerate a lost IRQ caused by IRQ_DISABLED

# 进程同步/通信

## A:SYSTEM V IPC:

### 1. Semaphore 信号量
int semctl(int sem_id, int sem_num, int command, ...);
int semget(key_t key, int num_sems, int sem_flags);
int semop(int sem_id, struct sembuf *sem_ops, size_t num_sem_ops);

### 2.Share Memory 共享内存
void *shmat(int shm_id, const void *shm_addr, int shmflg);
int shmctl(int shm_id, int cmd, struct shmid_ds *buf);
int shmdt(const void *shm_addr);
int shmget(key_t key, size_t size, int shmflg);

### 3.Message Queues 消息队列
int msgctl(int msqid, int cmd, struct msqid_ds *buf);
int msgget(key_t key, int msgflg);
int msgrcv(int msqid, void *msg_ptr, size_t msg_sz, long int msgtype, int msgflg);
int msgsnd(int msqid, const void *msg_ptr, size_t msg_sz, int msgflg);

## B:PIPE
## C:FIFO
## D:POSIX MSGQ
## E:SIGNAL

# 线程同步

### 1. semophore 信号量
sem_wait(&bin_sem);
sem_post(&bin_sem);
sem_destroy(&bin_sem);

### 2. Mutex 互斥
int pthread_mutex_init(pthread_mutex_t *mutex, const pthread_mutexattr_t *mutexattr);
int pthread_mutex_lock(pthread_mutex_t *mutex));
int pthread_mutex_unlock(pthread_mutex_t *mutex);
int pthread_mutex_destroy(pthread_mutex_t *mutex);

### 3. 读写锁
int pthread_rwlock_init(pthread_rwlock_t *restrict rwlock, const pthread_rwlockattr_t *restrict attr);
int pthread_rwlock_destroy(pthread_rwlock_t *rwlock);
int pthread_rwlock_rdlock(pthread_rwlock_t *rwlock);
int pthread_rwlock_wrlock(pthread_rwlock_t *rwlock);
int pthread_rwlock_unlock(pthread_rwlock_t *rwlock);
int pthread_rwlock_tryrdlock(pthread_rwlock_t *rwlock);
int pthread_rwlock_trywrlock(pthread_rwlock_t *rwlock);

### 4. 条件变量
pthread_cond_t cond = PTHREAD_COND_INITIALIZER;
int pthread_cond_init(pthread_cond_t *cond, pthread_condattr_t *cond_attr);
int pthread_cond_signal(pthread_cond_t *cond);
int pthread_cond_broadcast(pthread_cond_t *cond);
int pthread_cond_wait(pthread_cond_t *cond, pthread_mutex_t *mutex);
int pthread_cond_timedwait(pthread_cond_t *cond, pthread_mutex_t *mutex, const struct timespec *abstime);
int pthread_cond_destroy(pthread_cond_t *cond);
**Event**
**Critical section**

**自旋锁: 保护不同的CPU并发访问**
1, ISR; 2, forbid deferable function, 3, set positive preempt counter.
?只有当内核在执行异常处理程序而且抢占没有被禁止，才可能抢占内核?
**中断返回的时候会判断是返回到用户模式还是内核模式.在非抢占系统中只有返回用户模式才会判断TIF_NEED_RESCHED以及重新调度,返回内核态的时候不会(用户态抢占).而抢占式内核则不论中断返回用户态还是内核态都会判断并调度(内核态抢占)**
临界区 关中断,禁止抢占,
进入内核态的途径:
竞争条件的产生: 交叉内核控制路径; 多CPU

1,中断hdlr和tasklet不必可重入;2,仅被软中断和tasklet访问的每CPU变量不需要同步;3,仅被一种tasklet访问的数据结构不需要同步.

**内核同步:**
主要是防止多核处理器同时访问修改某段代码，或者在对设备驱动程序进行临界区保护。主要有一下几种方式：

1,每cpu变量,禁止抢占的情况下访问
2,原子操作,atomic_xxx(); xxx_bit()
3,内存屏障, 流水处理技术导致的指令交错.(优化屏障:避免编译器优化导致的指令重排序.(volatile))
4,自旋锁, (读-写spinlock,对读,写操作分别加锁)(顺序锁,写具有较高悠闲级)
内核路径自旋&进程挂起
5,信号量, (mutex是一种特殊的sem)(读写信号量,分别对读写设置)
6,顺序锁
7,本地中断禁止
8,本地软中断禁止
9,读-拷贝-更新, 读和拷贝同时进行,当读锁全部释放才进行更新.

**Logical Address:**
**Linear Address**: virtual address, 4G
**Physical Address**:
程序员可见Logical Address ->（分段单元）Linear Address ->（分页单元）Physical Address
Memory arbiter: DMA和CPU，CPU和CPU之间并发访问；由硬件访问，对程序员不可见
**Real mode**: 操作系统自举
**Protected mode**:
**段选择符**: 16 bit。Segment Selector(16 bit) + Offset(32 bit)。组成：
index, 指定放在GDT或者LDT中的相应**段描述符**入口。段描述符地址 = gdtr或ldtr + index * 8
TI, GDT还是LDT
RPL 特权级

**段寄存器**: 存放**段选择符**
CS, 当前特权级(Current Priviledge Level)使用2 bit，表明CPU当前特权。在Linux中，0为内核态，3为用户态。
SS, DS;
ES, FS, GS //都是附加段寄存器
SI, DI //源/目变址寄存器

**段描述符**Segment Descriptor, 8 byte (64 bit),
组成：描述段的特征（段线性地址，粒度，最后一个内存单元偏移，段类型(系统段或者普通代码/数据段)，类型和存取权限，
**段描述符特权DPL**，规定访问该段的特权级别
**CPL**, 当前进程的特权级别
Segment-Present标志，D/B标志(段偏移量32或者16位？)，AVL(Linux未用)）
**GDT** (Global Descriptor Table) 全局描述符表 gdtr控制寄存器, 最大8191 个(13 bit)
**LDT** (Local Descriptor Table) 局部描述符表 在不同的进程中创建, ldtr控制寄存器,
描述符种类：**代码段描述符，数据段描述符，任务状态段描述符TSSD**(该段用于保存处理器寄存器的内容，GDT中)，局部描述符表描述符LDTD(GDT 中)。
快速访问段描述符 缓存**段描述符**

**分段单元**：逻辑地址转换到线性地址 1: gdtr或ldtr + index * 8 寻址段描述符；2: 段描述符base字段+逻辑地址偏移
分段：不同进程分配不同的线性地址
分页：同一线性地址映射到不同的物理地址
用户代码段，用户数据段，内核代码段，内核数据段
宏__USER_CS, __USER_DS, __KERNEL_CS, __KERNEL_DS

**Linux GDT**
数组cpu_gdt_table 存放GDT
数组cpu_gdt_descr 存放GDT地址极其大小

**任务状态段TSS**: init_tss. TSS为246字节长
切换CPU时恢复现场；执行I/O时检查许可位图

缺省局部描述符表段 LDT
局部线程存储段 TLS 允许多线程应用使用最多3个局部于线程的数据段
**set_thread_area(), get_thread_area()**

**高级电源管理段3个**
PnP功能段5个
双重错误段 处理异常时引发另外一个异常，产生双重错误

**Linux LDT**
缺省LDT: default_ldt. 5个表项; modify_ldt()创建自己的局部描述符表

**调用门** 在调用预定义函数的时候改变CPU特权
**分页单元: page frame, page table**
**缺页异常** 访问类型与线性地址的访问权限不一致
控制寄存器cr0: PG标志为0时，线性地址被解释成物理地址
32位线性地址：Directory(10 bit) + Table(10 bit) + Offset(12 bit)

**页目录** page directory
组成 Present, Accessed, Dirty, R/W, User/Supervisor, PCD/PWT, Page size, Global
TLB: Translation Lookaside Buffer, Associative Memory, 俗称 快表

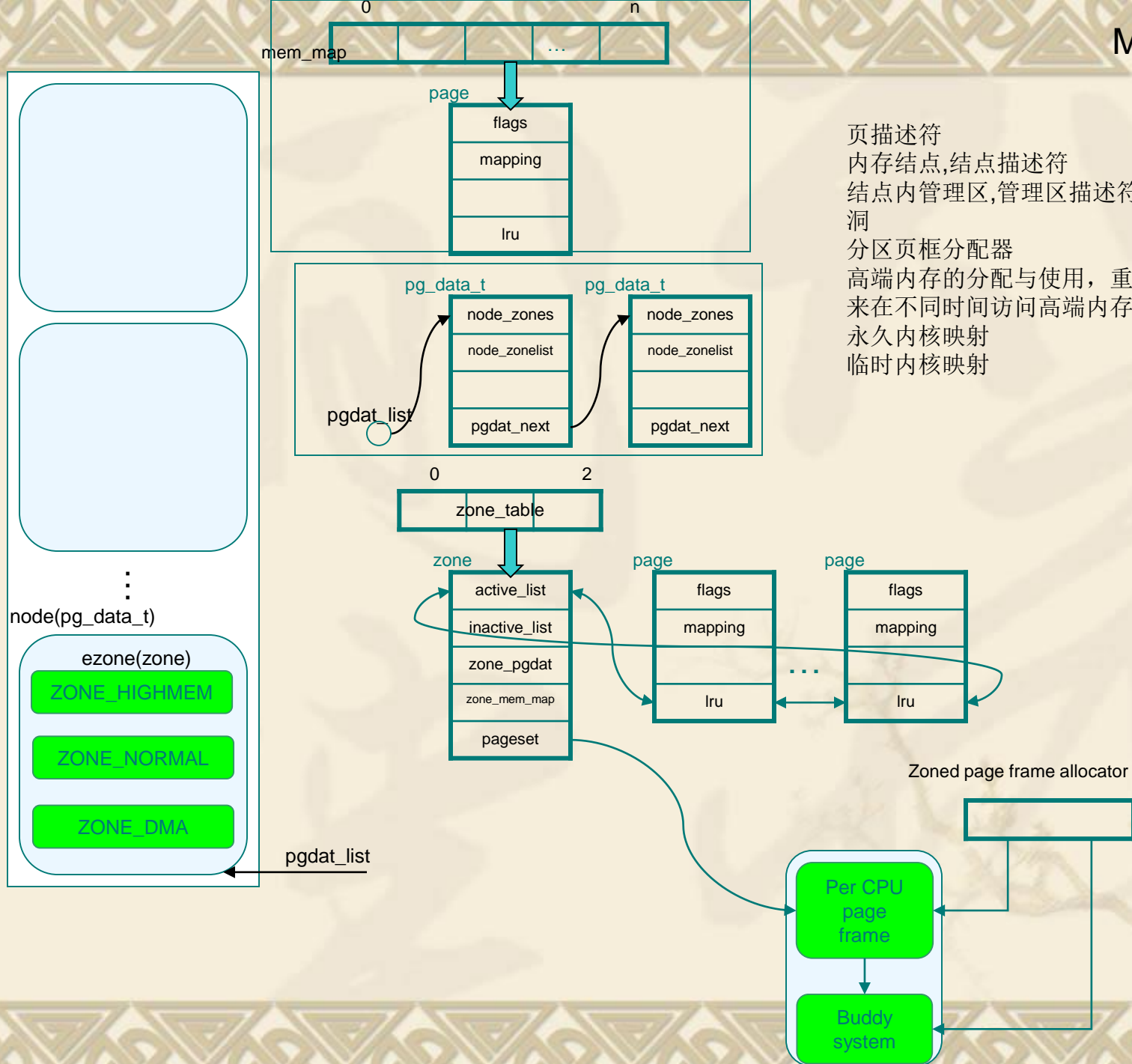**页表** page table
控制寄存器cr3: 正在使用页目录的物理地址
**扩展分页**
线性地址没有Table部分，页框大小为4M, page size置位，20位物理地址只有高10位
有意义。Cr4的PSE标志使得两种分页共存。
User/Supervisor=0时，只有内核态才能寻址。为1时，总能寻址
**物理地址扩展分页(PAE和PSE(Linux未使用))**
启用: cr4的PAE标志置位; page directory中的page size使用大尺寸页面.此时页面大小
为2M.

Memory

mem_map

page
flags
mapping
lru

pg_data_t
node_zones
node_zonelist
pgdat_next

pg_data_t
node_zones
node_zonelist
pgdat_next

pgdat_list

zone_table

zone
active_list
inactive_list
zone_pgdat
zone_mem_map
pageset

page
flags
mapping
lru

page
flags
mapping
lru

node(pg_data_t)

ezone(zone)
ZONE_HIGHMEM
ZONE_NORMAL
ZONE_DMA

pgdat_list

页描述符
内存结点,结点描述符
结点内管理区,管理区描述符
洞
分区页框分配器
高端内存的分配与使用，重复使用最后128M线性地址
来在不同时间访问高端内存(与高端内存映射)
永久内核映射
临时内核映射

Zoned page frame allocator

Per CPU page frame

Buddy system