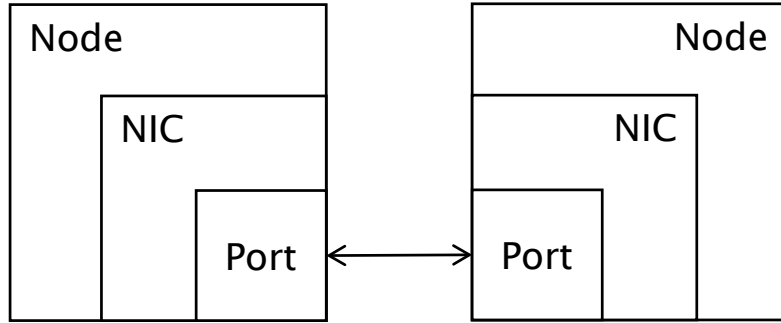


# Life of a packet

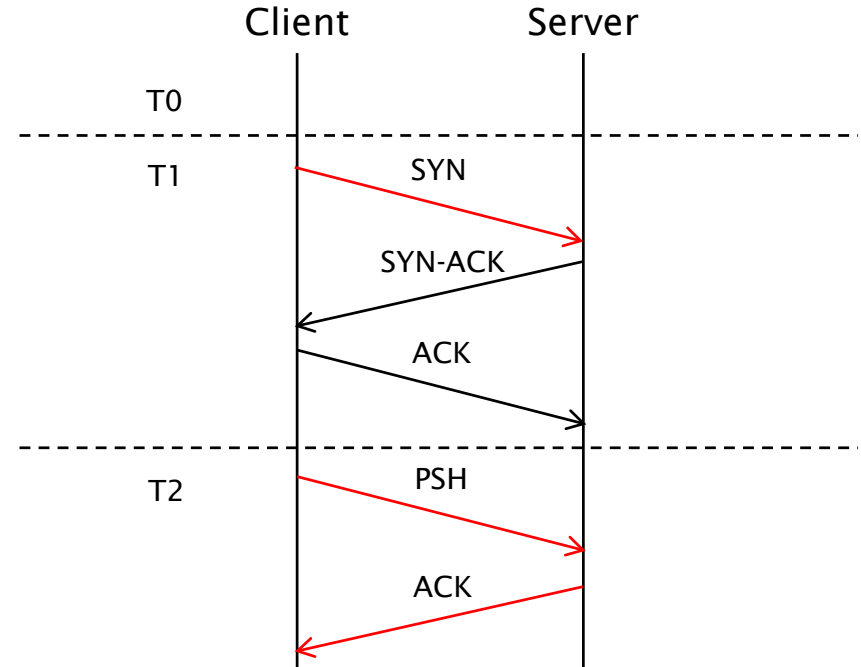
Xiao Jia

Pygmal Technologies

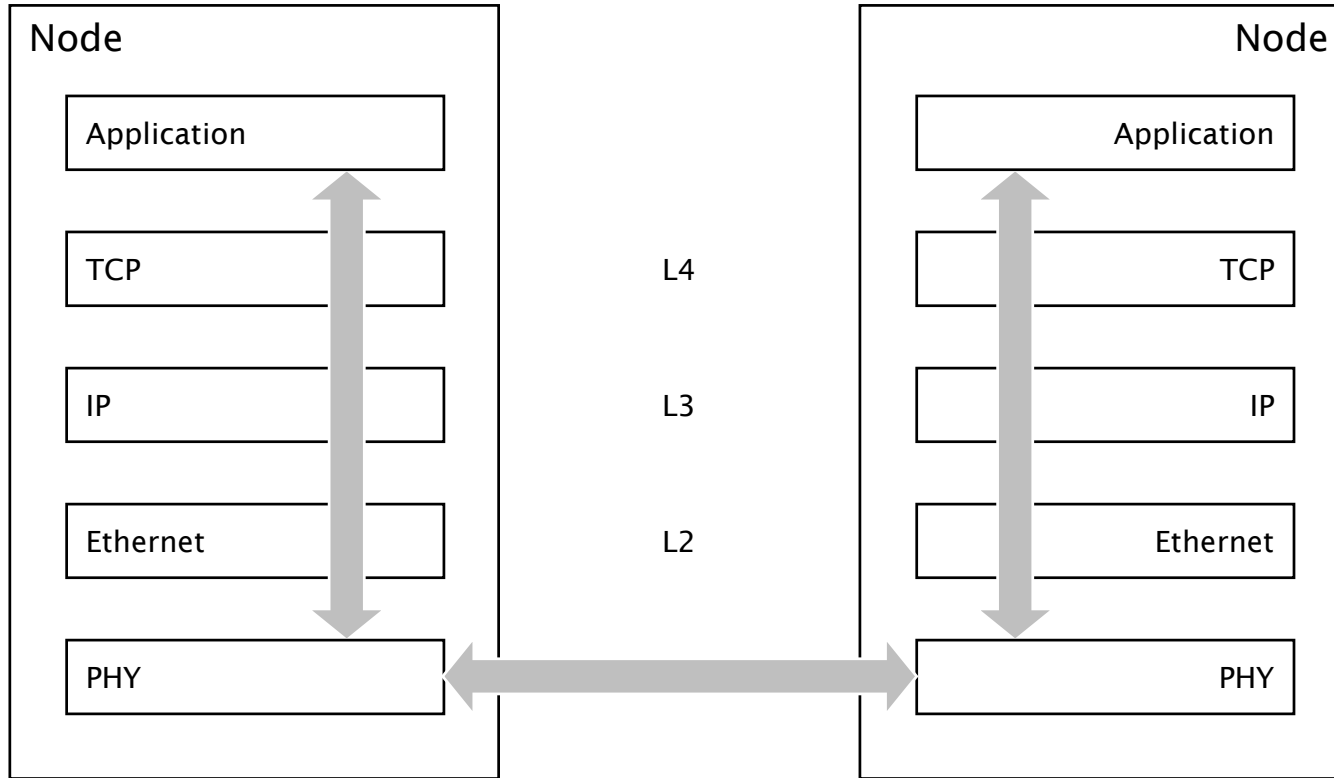
# Scope



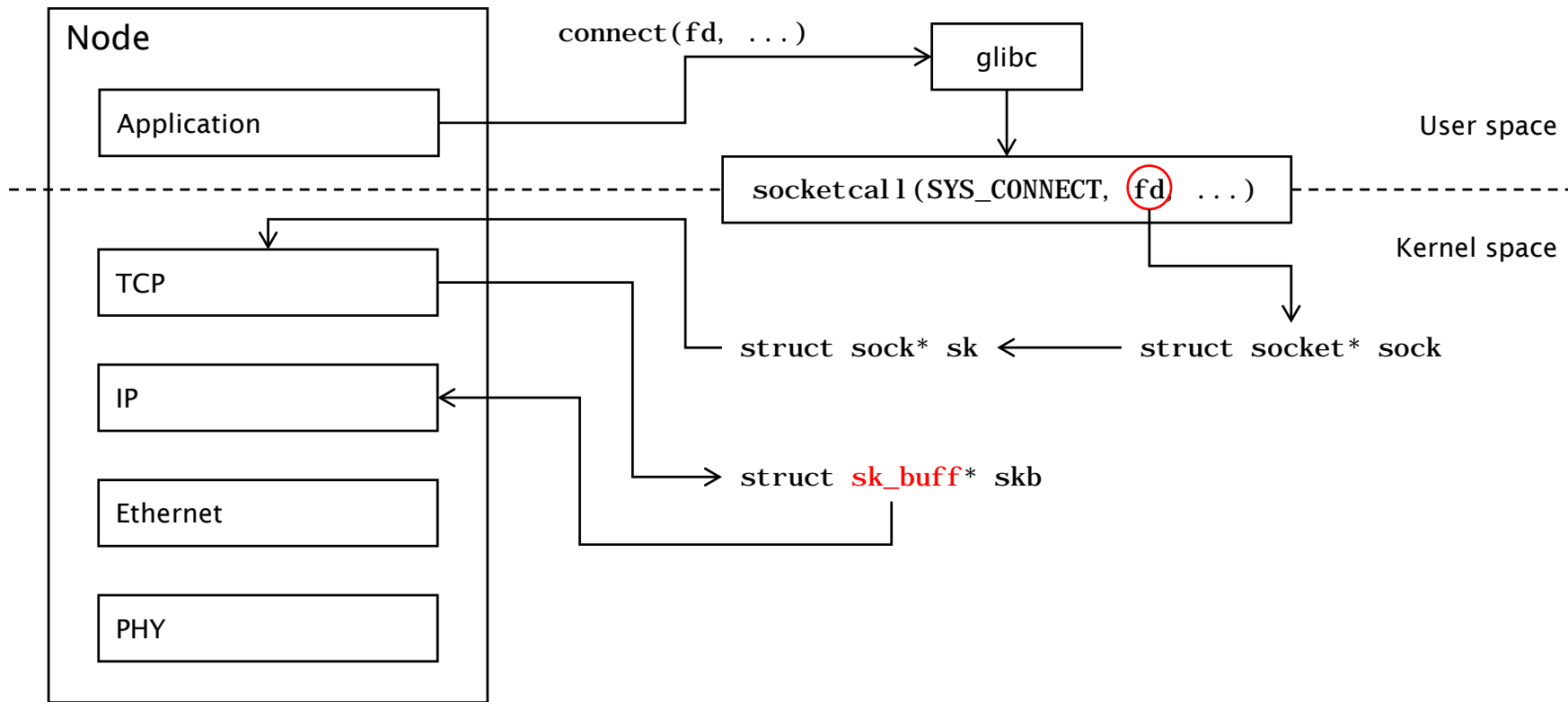
Time	Client	Server
T0		bind + listen
T1	connect	accept
T2	write	read



# Layers



# User/kernel space transition



# SKB (1)

```
struct sk_buff {
```

```
    struct sk_buff    *next;
```

```
    struct sk_buff    *prev;
```

```
    struct sk_buff_head *list;
```

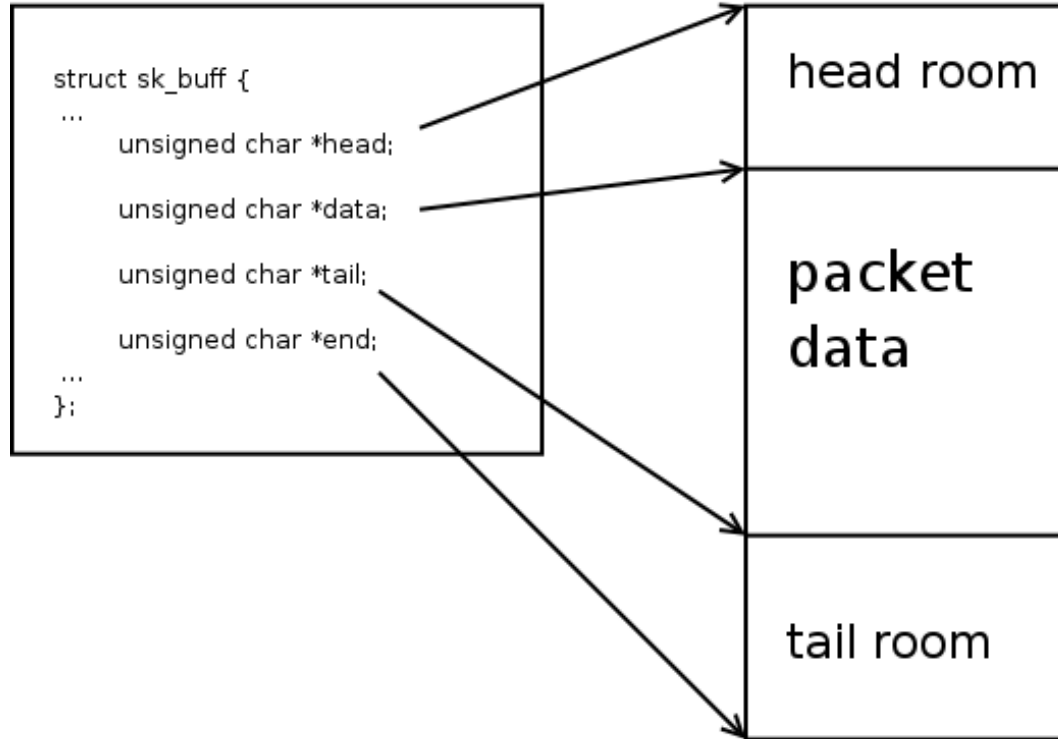
```
    struct sock        *sk;           for e.g. proper memory accounting
```

```
    struct net_device  *dev;
```

```
    struct dst_entry   *dst;
```

```
    char               cb[40];        Control block. TCP uses this to store  
                                     seq. number and retransmission state.
```

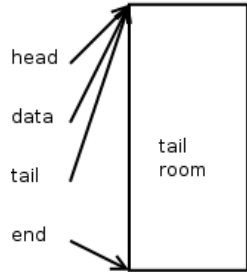
# SKB data area (1)



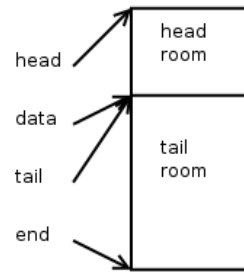
Source: [http://vger.kernel.org/~davem/skb\\_data.html](http://vger.kernel.org/~davem/skb_data.html)

# SKB data area (2)

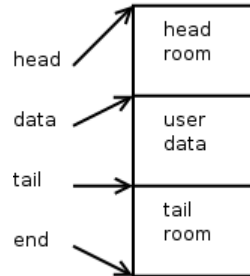
```
skb = alloc_skb(len, GFP_KERNEL);
```



```
skb_reserve(skb, header_len);
```

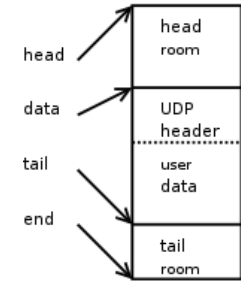


```
unsigned char *data = skb_put(skb, user_data_len);
int err = 0;
skb->csum = csum_and_copy_from_user(user_pointer, data,
                                     user_data_len, 0, &err);
if (err)
    goto user_fault;
```



```
struct inet_sock *inet = inet_sk(sk);
struct flowi *fl = &inet->cork.fl;
struct udphdr *uh;
```

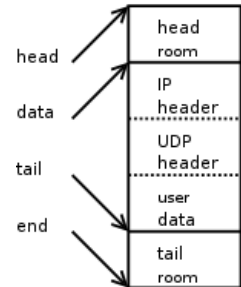
```
skb->h.raw = skb_push(skb, sizeof(struct udphdr));
uh = skb->h.uh;
uh->source = fl->fl_ip_sport;
uh->dest = fl->fl_ip_dport;
uh->len = htons(user_data_len);
uh->check = 0;
skb->csum = csum_partial((char *)uh,
                        sizeof(struct udphdr), skb->csum);
uh->check = csum_tcpudp_magic(fl->fl4_src, fl->fl4_dst,
                             user_data_len, IPPROTO_UDP, skb->csum);
if (uh->check == 0)
    uh->check = -1;
```



```
struct rtable *rt = inet->cork.rt;
struct iphdr *iph;
```

```
skb->nh.raw = skb_push(skb, sizeof(struct iphdr));
iph = skb->nh.iph;
iph->version = 4;
iph->ihl = 5;
iph->tos = inet->tos;
iph->tot_len = htons(skb->len);
iph->frag_off = 0;
iph->id = htons(inet->id++);
iph->ttl = ip_select_ttl(inet, &rt->u.dst);
iph->protocol = sk->sk_protocol; /* IPPROTO_UDP in this case */
iph->saddr = rt->rt_src;
iph->daddr = rt->rt_dst;
ip_send_check(iph);
```

```
skb->priority = sk->sk_priority;
skb->dst = dst_clone(&rt->u.dst);
```



## SKB (2)

```
struct sk_buff {
```

```
    union {
```

```
        struct tcphdr  *th;
```

```
        struct udphdr  *uh;
```

```
        struct icmphdr *icmph;
```

```
        struct igmp_hdr *igmp;
```

```
        struct iphdr   *iph;
```

```
        struct ipv6hdr *ipv6h;
```

```
        unsigned char  *raw;
```

```
    } h;
```

```
    union {
```

```
        struct iphdr   *iph;
```

```
        struct ipv6hdr *ipv6h;
```

```
        struct arphdr  *arph;
```

```
        unsigned char  *raw;
```

```
    } nh;
```

```
    union {
```

```
        unsigned char  *raw;
```

```
    } mac;
```



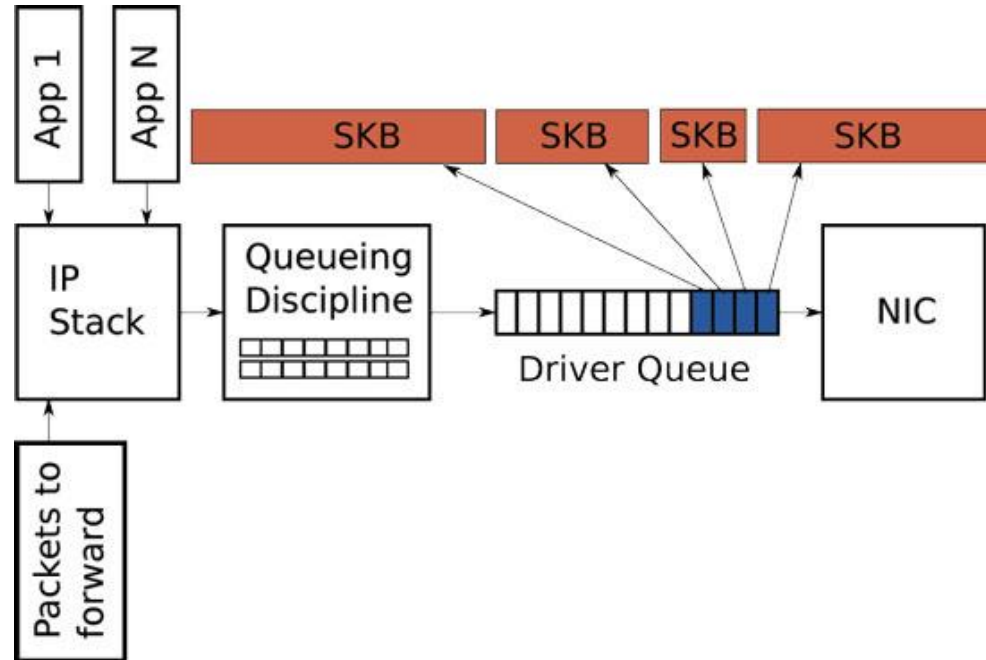
# Transmission path

- write/sendto/sendmsg
- tcp\_sendmsg
- tcp\_transmit\_skb
- ip\_queue\_xmit
- ip\_output
- dev\_queue\_xmit
- netif\_schedule
- dev\_hard\_start\_xmit
- dev\_kfree\_skb\_irq

# Receive path

- netif\_rx\_schedule
- net\_rx\_action
- netif\_receive\_skb
- ip\_rcv
- ip\_local\_deliver
- tcp\_v4\_rcv
- tcp\_rcv\_established
- skb\_copy\_datagram\_iovec
- read/recvfrom/recvmsg

# Queueing



Source: <http://www.linuxjournal.com/content/queueing-linux-network-stack>

# Debugging tools

- netperf
- google/neper
- perf
- tcpdump
- wireshark
- strace
- ethtool
- kprobe
- ftrace

# Missing pieces (a lot of them)

- qdisc
- netfilter (iptables)
- GRO
- RFS
- NAPI
- bonding
- Protocol details (e.g. TCP retransmission)
- ... ..