Netlink Design Document

by	Sean Ng	
mentor	Alexander Chernikov	
date	14th June 2021	

Description	
Overview of Design	
Design Details	3
1.Installing Netlink with netisr (used for deferred dispatch)	3
2.Registering Netlink as a socket (VNET)	4
3.Attaching the netlink socket	4
4.Sending a message	5
5.Netisr packet input	5
6.Netlink connect	6
7.Initialize netlink generic functionality	6
Design Decisions	7
Development Milestones	8
Phase 1	8
Phase 2	9
Phase 3	10
Phase 4	10

Description

Netlink is a linux kernel interface used for communication between userspace and kernel processes.

This project aims to port Netlink over to FreeBSD. The goal is to implement NETLINK_ROUTE in FreeBSD, and eventually support NETLINK_GENERIC as well. Eventually, the project aims to have a program use netlink. Currently, systems in FreeBSD which use Netlink sockets in their equivalent implementations in linux are currently supported by routing sockets.

Project is part of GSOC 2021. More details can be found in GSOC 2021

Overview of Design

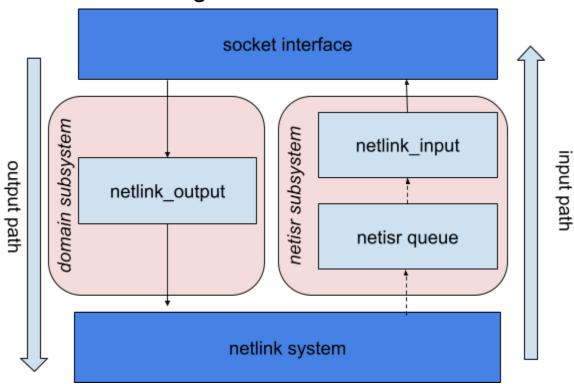


Figure: Overview of the various components needed to interface with the netlink system

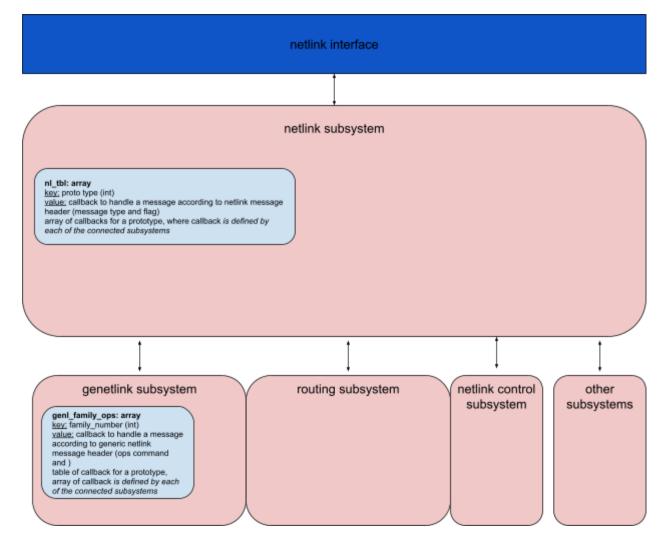


Figure: Overview of the netlink system

Design Details

1.Installing Netlink with netisr (used for deferred dispatch)

Description	netisr is the kernel network dispatch service. netisr_register is a kernel routine that manages a netisr_handler.
<u>Structs</u>	Implement netisr_handle: - nh_handler: field which is called for all input, refer to packet input (from kernel) section - nh_proto: new unique protocol number - nh_policy: use NETISR_POLICY_SOURCE

<u>Functions</u>	netlink_input: handles messages from kernel. Refer to netlink_input sysctl_netlink_netisr_maxqlen is used when checking queue, needed to register in netisr
man pages	netisr man pageSYSINIT is useful
TODO	

2.Registering Netlink as a socket (VNET)

Description	domain is an abstraction that network protocols are installed under (i.e. inetdomain/localdomain). each domain defines an array of protocol switch structures, one for each socket type.	
<u>Structs</u>	Implement domain - dom_family as PF_NETLINK - dom_protosw as struct defined in next bullet pointj Implement protosw - pr_output: function output to protocol, which we will define as netlink_output - pr_type: SOCK_RAW - pr_usrreqs: pr_usrreqs - refer to rt_sock for the others Implement pr_usrreqs Implement pr_usrreqs	
	static struct pr_usrreqs netlink_usrreqs = {	

```
can be equal to the process ID only for at most one socket
                     There are two ways to assign nl_pid to a netlink socket. If the application sets nl_pid before calling bind(2), then it is up to
                     the application to make sure that nl_pid is unique. If the
                     application sets it to 0, the kernel takes care of assigning it.
                     The kernel assigns the process ID to the first netlink socket the
                     process opens and assigns a unique nl_pid to every netlink socket
                     that the process subsequently creates.
                  from netlink man page
                 note that the length is updated in the process of the syscall
                      .pru detach = netlink detach,
                      .pru_disconnect = netlink_disconnect,
                      .pru peeraddr =
                                              netlink peeraddr,
                      .pru send = netlink send,
                      .pru shutdown = netlink shutdown,
                      .pru_sockaddr =
                                              netlink_sockaddr,
                      .pru_close =
                                            soisdisconnected
                  };
                  sending functions will be discussed in the later sections
man pages
                           domain man page
                           VNET_DOMAIN_SET/DOMAIN_SET is useful
TODO
                           What is exact difference between vnet vs non-vnet, or specifically, why
                           we call VNET DOMAIN SET instead of DOMAIN SET
                           protosw.ctloutput: I don't have a good idea on when this is called so I will
                           leave it out for now
                           Read https://flylib.com/books/en/2.849.1.145/1/
```

3. Attaching the netlink socket

<u>Description</u>	As with most socket APIs, upon attaching a socket, the raw socket for that address family needs to be instantiated and the fields filled up.	
<u>Structs</u>	Need to specify and implement a data structure to keep track of the protocols egistered for the netlink address family	
<u>Functions</u>	etlink_attach: retrieves raw socket pointer, and initializes the fields in the raw ocket	
man pages	https://flylib.com/books/en/2.849.1.150/1/	
TODO	need more documentation for dealing w raw sockets	

4.Sending a message

Description	User creates a message using mbuf, and sends it through netlink. When creating a packet for the socket, used both by protosw and netlinkisr		
Structs	Need to specify and implement a data structure to keep track of the protocols and callbacks for each netlink family		
Functions	netlink_output called with mbuf and socket 1. netlink receives packet, saves the protocol and the portid to the mbuf, 2. arrange packet content into a linear buffer 3. calling a callback on the linear buffer Note that both header and body are padded to the next 4 byte boundary Note that ACK messages need to be send for NLM_F_ACK messages Note that each message should have NLM_F_REQUEST		
man pages	https://datatracker.ietf.org/doc/html/rfc3549#section-2.3.2		
TODO	In Luigi's implementation, he: 1. copies out into a linear buffer, is it necessary for me to do this? 2. I don't understand his receive packet implementation. in netlink receive packet, each mbuf in the mbuf chain is assumed to have a packet header, and it appears that he calls the callback for EAC of these packets, which means that I should treat these packets as separate requests? Then why the while loop?		

5.Netisr packet input

Description	When processing an incoming packet from the netisr queue (nh_handler field in netisr_handler), function needs to implement callback that identifies the right port	
<u>Structs</u>	Implement sockaddr within the implementation of the nh_handler itself - From kernel, the pid of the sockaddress is 0. Implement sockproto_within the implementation of the nh_handler itself	
Functions	<pre>raw_input_netlink: handler defined in netisr_handler. Will create nl_src and nl_proto, using both to pass into raw_input_ext for the dispatcher to route the request properly raw_input_ext will map the sockaddr with the right port.</pre>	
man pages	sys/net/raw_cb.h sys/net/ <i>raw_usrreq.c</i> mbuf manpage (mtod is useful)	
TODO	Why do we bother with sockproto and not read the values directly?	

In Luigi, implementation he doesn't even use sockproto in
raw_input_netlink_cb?

6.Netlink connect

Description	implementation for a connect syscall	
<u>Structs</u>	Modify fields in the <i>socket</i> struct - nl_src_portid corresponds to so_fibnum - nl_dst_portid corresponds to so_user_cookie	
<u>Functions</u>	tlink_connect: sets the portid on the socket, call soisconnected, checks size	
man pages	sys/sys/socketvar.h	
TODO	Figure out how linux netlink handles the portnumbers	

7.Initialize netlink generic functionality

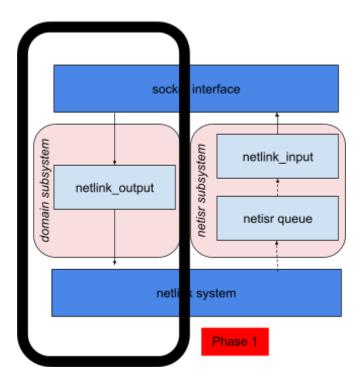
Description	Register new proto in the netlink family, initialize structs for generic. Will be similar to linux's implementation for uniformity.		
<u>Structs</u>	Define and implement a struct <i>genl_ops</i> that specifies registered generic message handlers: - u8 generic message type - doit and dumpit calllbacks - nla policy - more defined in <i>Theory and Implementation</i> book		
Functions	genetlinkload that initializes all the functions		
man pages	 LIST_INIT and related macros are useful SYS_INIT is useful https://elixir.bootlin.com/linux/latest/source/include/net/genetlink.h#L48 Theory and Implementation book page 27 		
TODO	figure out exact difference between the doit callback and the dumpit callback		

Design Decisions

- 1. Decision: Linux uses an <u>netlink_sock struct</u> to store details relevant to a connection such as dst_port_id, but this pattern does not appear to be a common pattern in FreeBSD sockets, hence I will not use this pattern. Currently, these values are only needed once and hence I will parse the header to retrieve these values when needed, which is the same as Luigi's implementation.
- 2. Decision: Data structure for generic netlink callbacks. Luigi uses <u>a single list</u> that he iterates on to find the right generic family structure for a generic generic family id/name. Instead, I decided to use an array of families, similar to that of the <u>structure that stores</u> netlink callbacks. Linux uses a radix tree (linux implementation) (this is an idr)

Development Milestones

Phase 1



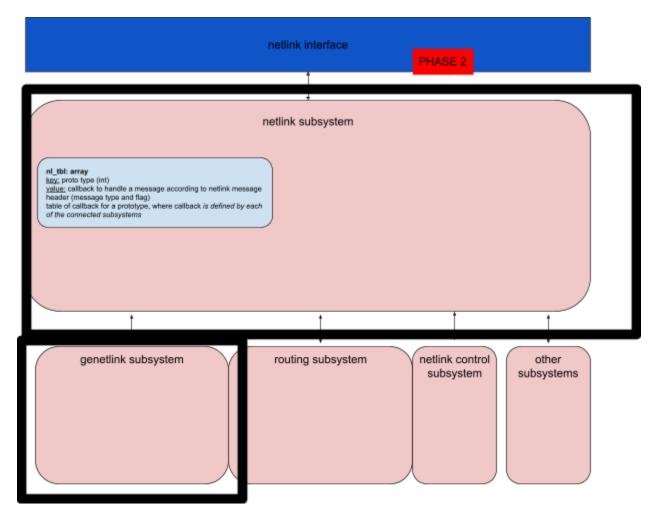
Description: Define interfaces to interact the domain subsystem with the netlink subsystem

Goal: Calls the functions when the sockets are open and message is sent

Time: 1 week

Detailed Description Section 2, 3, 4

Phase 2



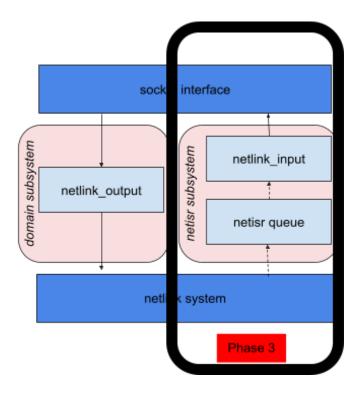
Description: Defining the netlink subsystem, defining the genetlink subsystem and interfacing both. Does not have to pimplement detailed control subsystem for generic netlink.

Goal: Ensure that the right functions are called when calling a generic socket

Time: 1-2 weeks

Detailed Description Section 4

Phase 3



Description: Interfacing with netisr subsystem

Goal: Ensures that the socket interface should receive a message

Time: 1 week

Detailed Description Section 1, 5

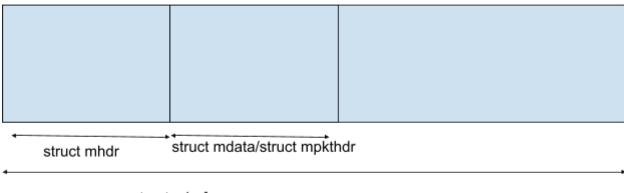
Phase 4

[TBC, will confirm at a later date]

Appendix

My understanding about each packet

mbuf that is not extended, with pkthdr



struct mbuf

1 mbuf chain, 1 packet

Previous header (to be phased out)

```
A message header consists of one of the following:
.Bd -literal
struct rt_msghdr {
    u_short rtm_msglen;
                               /* to skip over non-understood messages */
    u_char rtm_version;
                               /* future binary compatibility */
    u_char rtm_type;
                               /* message type */
    u_short rtm_index;
                               /* flags, incl. kern & message, e.g. DONE */
           rtm_flags;
                               /* bitmask identifying sockaddrs in msg */
           rtm_addrs;
    pid_t
           rtm_pid;
                               /* for sender to identify action */
           rtm_seq;
           rtm_errno;
                               /* why failed */
           rtm_fmask;
                               /* bitmask used in RTM_CHANGE message */
                               /* which metrics we are initializing */
    u_long rtm_inits;
    struct rt_metrics rtm_rmx;> /* metrics themselves */
struct if_msghdr {
```

/* rtm type */

in route.h

/* rtm flags */

```
/* TTM__flags */

172 #define RTF_UP 0x1 /* route usable */
173 #define RTF_GATEWAY 0x2 /* destination is a gateway */
174 #define RTF_GATEWAY 0x2 /* destination is a gateway */
175 #define RTF_DEVA 0x4 /* host entry (net otherwise) */
175 #define RTF_DYNAMIC 0x10 /* created dynamically (by redirect) */
176 #define RTF_DNE 0x20 /* modified dynamically (by redirect) */
177 #define RTF_DONE 0x40 /* message confirmed */
178 /* 0x80 unused, was RTF_DELCLONISG */
180 /* 0x100 unused, was RTF_DELCLONISG */
181 #define RTF_XRESOLVE 0x200 /* external daemon resolves name */
182 #define RTF_LLINFO 0x400 /* DEPRECATED - exists ONLY for backward
183 compatibility */
184 #define RTF_LLINFO 0x400 /* DEPRECATED - exists ONLY for backward
185 #define RTF_LLINFO 0x400 /* manually added */
186 #define RTF_BLACKHOLE 0x100 /* just discard pkts (during updates) */
187 #define RTF_PROTO2 0x4000 /* protocol specific routing flag */
189 /* 0x10000 unused, was RTF_PRCLONING */
191 #define RTF_PROTO3 0x40000 /* protocol specific routing flag */
192 #define RTF_PROTO3 0x40000 /* route is minumulable */
193 #define RTF_PROTOA 0x400000 /* route represents a local address */
196 #define RTF_BLACKLOAST 0x4000000 /* route represents a boast address */
196 #define RTF_BROTOACAST 0x400000 /* route represents a boast address */
197 /* 0x8000000 and up unassigned */
198 #define RTF_STICKY 0x10000000 /* always route dst->src */
199
0000 0x4000000 unused, was RTF_RNH_LOCKED */
   199
200 /*
201
                                               0x40000000 unused, was RTF_RNH_LOCKED */
   202 #define RTF_GWFLAG_COMPAT 0x80000000 /* a compatibility bit for interacting
                                                                            with existing routing apps *
     205 /* Mask of RTF flags that are allowed to be modified by RTM_CHANGE. */
     206 #define RTF_FMASK \
207 (RTF_PROTO1 | RTF_PROTO2 | RTF_PROTO3 | RTF_BLACKHOLE | \
                              RTF_REJECT | RTF_STATIC | RTF_STICKY)
```

/*rtm addrs*/

/*rtm inits*/

currently using

```
28 struct rtmsg {
                           rtm_family;
       unsigned char
       unsigned char
                           rtm_dst_len;
       unsigned char
                           rtm_src_len;
       unsigned char
                           rtm_tos;
       unsigned char
                           rtm_table; /* Routing table id */
       unsigned char
                           rtm_protocol; /* Routing protocol; see below */
                           rtm_scope; /* See below */
       unsigned char
                                       /* See below */
       unsigned char
                           rtm_type;
238
       unsigned
                       rtm_flags;
240 };
```

rtm family

```
rtm-><mark>rtm_family</mark> = AF_INET;
rtm->rtm_family = AF_INET6;
```

nlmsg_type

```
22 /* Types of messages */
23 /* Types of messages */
23 /* Types of messages */
25 /* RTM_BASE = 16,
26 //#define RTM_BASE RTM_BASE
27 //
28 // RTM_NEWLINK = 16,
29 //#define RTM_MENLINK RTM_NEWLINK
30 // RTM_DELLINK,
31 //#define RTM_OELLINK RTM_DELLINK
32 // RTM_GETLINK,
32 //#define RTM_GETLINK
34 // RTM_SETLINK
35 //#define RTM_SETLINK
36 // RTM_NEWADDR = 20,
38 //#define RTM_NEWADDR RTM_NEWADDR
39 // RTM_DELADDR,
40 //#define RTM_DELADDR
41 // RTM_GETADDR,
41 // RTM_GETADDR,
42 //#define RTM_DELADDR RTM_DELADDR
41 // RTM_GETADDR,
41 // RTM_GETADDR,
42 //#define RTM_DELADDR RTM_DELADDR
```

rtm_type

```
enum {
    RTN_UNSPEC,
    RTN_UNICAST,
    RTN_ENCAST,
    RTN_ENCAST,
    RTN_ENCAST,
    RTN_BROADCAST,
    RTN_BROADCAST,
    RTN_ANYCAST,
    RTN_ENCAST,
    RTN_ENCAST,
```

rtm_protocol

rtm_flags

```
#define RTM F NOTIFY
                            0x100 /* Notify user of route change */
#define RTM F CLONED
                             0x200 /* This route is cloned
#define RTM_F_EQUALIZE
                             0x400 /* Multipath equalizer: NI */
#define RTM F PREFIX
                            0x800 /* Prefix addresses
#define RTM F LOOKUP TABLE 0x1000 /* set rtm_table to FIB lookup result */
#define RTM_F_FIB_MATCH
                                0x2000 /* return full fib lookup match */
#define RTM F OFFLOAD
                             0x4000 /* route is offloaded */
#define RTM_F_TRAP
                        0x8000 /* route is trapping packets */
#define RTM F OFFLOAD FAILED 0x20000000 /* route offload failed, this value
              * is chosen to avoid conflicts with
              * other flags defined in
              * include/uapi/linux/ipv6 route.h
```

need to convert to

sockaddr array

Field (rt_addr_info)	Source (rt_msghdr)	New Source (rtmsg)
rti_mflags	rtm_inits	
rtm_rmx	rtm_rmx	nlattr
rti_flags	rtm_flags	rtm_flags [or can't convert?]
rti_addrs	rtm_addrs	
rti_info	socket addresses are layed out in order after the rtmsg_hdr. it is "initialized" by considering the flags, and will be in the same order	populated using nl attributes

kgdb /boot/kernel/kernel /var/crash/vmcore.last ./rtnl-route-add xn0 10.0.1.12 32 172.31.18.62 ./rtnl-route-get 10.0.1.12 route delete 10.0.1.12

Source	Source (rt_msghdr)	New Source (rtmsg)	
	rtm_type	nlmsg_type	
	rtm_falgs (subset)	rtm_type	

Table for RTM GETROUTE

Problem: rtsock uses rib_cmd_info and nhop_object to update rt_msghdr

I need to use: rib_cmd_info and nh_object to update struct rtmsg

Source(rt_msghdr)	Linux Source (struct fib_rt_info)	Output (rtmsg)
	fri->type	rtm->rtm_type
	fri->scope	rtm->rtm_scope
	fri->fi->fib_flags	rtm->rtm_flags
look at rib_action	fri->tb_id	rtm->rtm_table
	fi->rtm_protocol	rtm->rtm_protocol

What corresponds to the table?

```
328 struct sockaddr {
329 unsigned char sa_len; /* total length */
330 sa_family_t sa_family; /* address family */
331 char sa_data[14]; /* actually longer; address value */
332 };
```

Hey Alex,

I got something simple for RTM_GETROUTE working!

1. Note: There does not appear to be corresponding mappings for the following fields in "struct rtmsg" (netlink's routing message).

```
rtm_type - linux retrieves this straight from its version of const struct fib_rt_info. freebsd's fib_rt_info does not have this information rtm_table - refers to routing table, and I don't know where this information is. rtm_scope - the fields don't make sense to freebsd rtm_flags - [not the same as rt_msg_hdr's rtm_flags] nothing corresponding
```

2. Is there one routing table? Or multiple routing table?

```
OIF = rtm index
IFA/IFP not present in rtmsg
is PRIORITY = Weight?
 33
      if (tb[RTA TABLE]) {
 34
         printf("table=%u ", mnl_attr_get_u32(tb[RTA_TABLE]));
 35
      }
Not sure. Only 1 routing table?
      if (tb[RTA_DST]) {
 36
 37
         struct in_addr *addr = mnl_attr_get_payload(tb[RTA_DST]);
 38
         printf("dst=%s ", inet_ntoa(*addr));
 39
      }
Done
 40
      if (tb[RTA_SRC]) {
 41
         struct in_addr *addr = mnl_attr_get_payload(tb[RTA_SRC]);
 42
         printf("src=%s ", inet_ntoa(*addr));
 43
      }
 44
      if (tb[RTA OIF]) {
 45
         printf("oif=%u ", mnl_attr_get_u32(tb[RTA_OIF]));
 46
      }
```

Done

```
47
      if (tb[RTA FLOW]) {
 48
         printf("flow=%u ", mnl_attr_get_u32(tb[RTA_FLOW]));
 49
      }
Not sure
 50
      if (tb[RTA_PREFSRC]) {
 51
         struct in addr *addr = mnl attr get payload(tb[RTA PREFSRC]);
 52
         printf("prefsrc=%s ", inet_ntoa(*addr));
 53
      }
Not sure
 54
      if (tb[RTA GATEWAY]) {
 55
         struct in_addr *addr = mnl_attr_get_payload(tb[RTA_GATEWAY]);
 56
         printf("gw=%s ", inet_ntoa(*addr));
 57
      }
Done
      if (tb[RTA PRIORITY]) {
 58
 59
         printf("prio=%u ", mnl_attr_get_u32(tb[RTA_PRIORITY]));
 60
      }
Not sure
 61
      if (tb[RTA_METRICS]) {
 62
         int i;
 63
         struct nlattr *tbx[RTAX_MAX+1] = {};
 64
 65
         mnl_attr_parse_nested(tb[RTA_METRICS], data_attr_cb2, tbx);
 66
 67
         for (i=0; i<RTAX_MAX; i++) {
 68
           if (tbx[i]) {
 69
             printf("metrics[%d]=%u ",
 70
                i, mnl_attr_get_u32(tbx[i]));
 71
           }
 72
         }
 73
      }
 74 }
```