```
**********************************************************
   27817 Wed Sep 26 12:51:33 2012
new/usr/src/common/iscsit/iscsit_common.c
inet_pton
**********************************************************
   1 /*
   2  * CDDL HEADER START
   3  *
   4  * The contents of this file are subject to the terms of the
   5  * Common Development and Distribution License (the "License").
   6  * You may not use this file except in compliance with the License.
   7  *
   8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
   9  * or http://www.opensolaris.org/os/licensing.
  10  * See the License for the specific language governing permissions
  11  * and limitations under the License.
  12  *
  13  * When distributing Covered Code, include this CDDL HEADER in each
  14  * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
  15  * If applicable, add the following below this CDDL HEADER, with the
  16  * fields enclosed by brackets "[]" replaced with your own identifying
  17  * information: Portions Copyright [yyyy] [name of copyright owner]
  18  *
  19  * CDDL HEADER END
  20  */
  21 /*
  22  * Copyright (c) 2008, 2010, Oracle and/or its affiliates. All rights reserved.
  23  */
  24 /*
  25  * Copyright (c) 2012, Nexenta Systems, Inc. All rights reserved.
  26  */
  27 #endif /* ! codereview */

  29 #include <sys/time.h>

  31 #if defined(_KERNEL)
  32 #include <sys/ddi.h>
  33 #include <sys/types.h>
  34 #include <sys/sunddi.h>
  35 #include <sys/socket.h>
  36 #include <inet/tcp.h>
  37 #include <inet/ip.h>
  38 #endif /* ! codereview */
  39 #else
  40 #include <stdio.h>
  41 #include <strings.h>
  42 #include <stdlib.h>
  43 #include <errno.h>
  44 #include <sys/types.h>
  45 #include <sys/socket.h>
  46 #include <netinet/in.h>
  47 #include <arpa/inet.h>
  48 #endif

  50 #include <sys/iscsit/iscsit_common.h>
  51 #include <sys/iscsi_protocol.h>
  52 #include <sys/iscsit/isns_protocol.h>

  54 void *
  55 iscsit_zalloc(size_t size)
  56 {
  57 #if defined(_KERNEL)
  58         return (kmem_zalloc(size, KM_SLEEP));
  59 #else
  60         return (calloc(1, size));
  61 #endif
```

```
  62 }

  64 void
  65 iscsit_free(void *buf, size_t size)      /* ARGSUSED */
  66 {
  67 #if defined(_KERNEL)
  68         kmem_free(buf, size);
  69 #else
  70         free(buf);
  71 #endif
  72 }

  74 /*
  75  * default_port should be the port to be used, if not specified
  76  * as part of the supplied string 'arg'.
  77  */

  79 #define NI_MAXHOST      1025
  80 #define NI_MAXSERV      32


  83 struct sockaddr_storage *
  84 it_common_convert_sa(char *arg, struct sockaddr_storage *buf,
  85     uint32_t default_port)
  86 {
  87         /* Why does addrbuf need to be this big!??! XXX */
  88         char            addrbuf[NI_MAXHOST + NI_MAXSERV + 1];
  89         char            *addr_str;
  90         char            *port_str;
  91 #ifndef _KERNEL
  92         char            *errchr;
  93 #endif
  94         long            tmp_port = 0;
  95         sa_family_t     af;

  97         struct sockaddr_in      *sin;
  98         struct sockaddr_in6     *sin6;
  99         struct sockaddr_storage *sa = buf;

 101         if (!arg || !buf) {
 102                 return (NULL);
 103         }

 105         bzero(buf, sizeof (struct sockaddr_storage));

 107         /* don't modify the passed-in string */
 108         (void) strlcpy(addrbuf, arg, sizeof (addrbuf));

 110         addr_str = addrbuf;

 112         if (*addr_str == '[') {
 113                 /*
 114                  * An IPv6 address must be inside square brackets
 115                  */
 116                 port_str = strchr(addr_str, ']');
 117                 if (!port_str) {
 118                         /* No closing bracket */
 119                         return (NULL);
 120                 }

 122                 /* strip off the square brackets so we can convert */
 123                 addr_str++;
 124                 *port_str = '\0';
 125                 port_str++;

 127                 if (*port_str == ':') {
```

```
128                             /* TCP port to follow */
129                             port_str++;
130                     } else if (*port_str == '\0') {
131                             /* No port specified */
132                             port_str = NULL;
133                     } else {
134                             /* malformed */
135                             return (NULL);
136                     }
137                     af = AF_INET6;
138             } else {
139                     port_str = strchr(addr_str, ':');
140                     if (port_str) {
141                             *port_str = '\0';
142                             port_str++;
143                     }
144                     af = AF_INET;
145             }

147             if (port_str) {
148     #if defined(_KERNEL)
149                     if (ddi_strtol(port_str, NULL, 10, &tmp_port) != 0) {
150                             return (NULL);
151                     }
152     #else
153                     tmp_port = strtol(port_str, &errchr, 10);
154     #endif
155                     if (tmp_port < 0 || tmp_port > 65535) {
156                             return (NULL);
157                     }
158             } else {
159                     tmp_port = default_port;
160             }

162             sa->ss_family = af;

164             sin = (struct sockaddr_in *)sa;
165             if (af == AF_INET) {
166                     if (inet_pton(af, addr_str,
167                         (void *)&(sin->sin_addr.s_addr)) != 1) {
168                             return (NULL);
169                     }
 24                     /*
 25                      * intet_pton does not seem to convert to network
 26                      * order in kernel. This is a workaround until the
 27                      * inet_pton works or we have our own inet_pton function.
 28                      */
 29     #ifdef _KERNEL
 30                     sin->sin_addr.s_addr = ntohl((uint32_t)sin->sin_addr.s_addr);
 31     #endif
170                     sin->sin_port = htons(tmp_port);
171             } else {
172                     sin6 = (struct sockaddr_in6 *)sa;
173                     if (inet_pton(af, addr_str,
174                         (void *)&(sin6->sin6_addr.s6_addr)) != 1) {
175                             return (NULL);
176                     }
177                     sin6->sin6_port = htons(tmp_port);
178             }

180             /* successful */
181             return (sa);
182 }
_____unchanged_portion_omitted_
```

```
**********************************************************
   49721 Wed Sep 26 12:51:35 2012
new/usr/src/uts/common/fs/nfs/nfs4_srv_deleg.c
inet_pton
**********************************************************
   1 /*
   2  * CDDL HEADER START
   3  *
   4  * The contents of this file are subject to the terms of the
   5  * Common Development and Distribution License (the "License").
   6  * You may not use this file except in compliance with the License.
   7  *
   8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
   9  * or http://www.opensolaris.org/os/licensing.
  10  * See the License for the specific language governing permissions
  11  * and limitations under the License.
  12  *
  13  * When distributing Covered Code, include this CDDL HEADER in each
  14  * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
  15  * If applicable, add the following below this CDDL HEADER, with the
  16  * fields enclosed by brackets "[]" replaced with your own identifying
  17  * information: Portions Copyright [yyyy] [name of copyright owner]
  18  *
  19  * CDDL HEADER END
  20  */
  21 /*
  22  * Copyright 2009 Sun Microsystems, Inc.  All rights reserved.
  23  * Use is subject to license terms.
  24  */
  25 /*
  26  * Copyright 2012 Nexenta Systems, Inc. All rights reserved.
  27  */
  28 #endif /* ! codereview */

  30 #include <sys/systm.h>
  31 #include <rpc/auth.h>
  32 #include <rpc/clnt.h>
  33 #include <nfs/nfs4_kprot.h>
  34 #include <nfs/nfs4.h>
  35 #include <nfs/lm.h>
  36 #include <sys/cmn_err.h>
  37 #include <sys/disp.h>
  38 #include <sys/sdt.h>

  40 #include <sys/pathname.h>

  42 #include <sys/strsubr.h>
  43 #include <sys/ddi.h>

  45 #include <sys/vnode.h>
  46 #include <sys/sdt.h>
  47 #include <inet/common.h>
  48 #include <inet/ip.h>
  49 #include <inet/ip6.h>

  51 #define MAX_READ_DELEGATIONS 5

  53 krwlock_t rfs4_deleg_policy_lock;
  54 srv_deleg_policy_t rfs4_deleg_policy = SRV_NEVER_DELEGATE;
  55 static int rfs4_deleg_wlp = 5;
  56 kmutex_t rfs4_deleg_lock;
  57 static int rfs4_deleg_disabled;
  58 static int rfs4_max_setup_cb_tries = 5;

  60 #ifdef DEBUG
```

```
  62 static int rfs4_test_cbgetattr_fail = 0;
  63 int rfs4_cb_null;
  64 int rfs4_cb_debug;
  65 int rfs4_deleg_debug;

  67 #endif

  69 static void rfs4_recall_file(rfs4_file_t *,
  70     void (*recall)(rfs4_deleg_state_t *, bool_t),
  71     bool_t, rfs4_client_t *);
  72 static  void            rfs4_revoke_file(rfs4_file_t *);
  73 static  void            rfs4_cb_chflush(rfs4_cbinfo_t *);
  74 static  CLIENT          *rfs4_cb_getch(rfs4_cbinfo_t *);
  75 static  void            rfs4_cb_freech(rfs4_cbinfo_t *, CLIENT *, bool_t);
  76 static rfs4_deleg_state_t *rfs4_deleg_state(rfs4_state_t *,
  77     open_delegation_type4, int *);

  79 /*
  80  * Convert a universal address to an transport specific
  81  * address using inet_pton.
  82  */
  83 static int
  84 uaddr2sockaddr(int af, char *ua, void *ap, in_port_t *pp)
  85 {
  86         int dots = 0, i, j, len, k;
  87         unsigned char c;
  88         in_port_t port = 0;

  90         len = strlen(ua);

  92         for (i = len-1; i >= 0; i--) {

  94                 if (ua[i] == '.')
  95                         dots++;

  97                 if (dots == 2) {

  99                         ua[i] = '\0';
 100                         /*
 101                          * We use k to remember were to stick '.' back, since
 102                          * ua was kmem_allocateded from the pool len+1.
 103                          */
 104                         k = i;
 105                         if (inet_pton(af, ua, ap) == 1) {

 107                                 c = 0;

 109                                 for (j = i+1; j < len; j++) {
 110                                         if (ua[j] == '.') {
 111                                                 port = c << 8;
 112                                                 c = 0;
 113                                         } else if (ua[j] >= '0' &&
 114                                             ua[j] <= '9') {
 115                                                 c *= 10;
 116                                                 c += ua[j] - '0';
 117                                         } else {
 118                                                 ua[k] = '.';
 119                                                 return (EINVAL);
 120                                         }
 121                                 }
 122                                 port += c;

  26                                 /* reset to network order */
  27                                 if (af == AF_INET) {
  28                                         *(uint32_t *)ap =
```

**new/usr/src/uts/common/fs/nfs/nfs4_srv_deleg.c** 3

```
  29                                        htonl(*(uint32_t *)ap);
 124                                *pp = htons(port);
  31                        } else {
  32                                int ix;
  33                                uint16_t *sap;

  35                                for (sap = ap, ix = 0; ix <
  36                                    sizeof (struct in6_addr) /
  37                                    sizeof (uint16_t); ix++)
  38                                        sap[ix] = htons(sap[ix]);

  40                                *pp = htons(port);
  41                        }
 126                        ua[k] = '.';
 127                        return (0);
 128                } else {
 129                        ua[k] = '.';
 130                        return (EINVAL);
 131                }
 132            }
 133        }

 135        return (EINVAL);
 136 }
_____unchanged_portion_omitted_
```

```
**********************************************************
   24552 Wed Sep 26 12:51:38 2012
new/usr/src/uts/common/fs/sockfs/nl7c.c
inet_pton
**********************************************************
   1 /*
   2  * CDDL HEADER START
   3  *
   4  * The contents of this file are subject to the terms of the
   5  * Common Development and Distribution License (the "License").
   6  * You may not use this file except in compliance with the License.
   7  *
   8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
   9  * or http://www.opensolaris.org/os/licensing.
  10  * See the License for the specific language governing permissions
  11  * and limitations under the License.
  12  *
  13  * When distributing Covered Code, include this CDDL HEADER in each
  14  * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
  15  * If applicable, add the following below this CDDL HEADER, with the
  16  * fields enclosed by brackets "[]" replaced with your own identifying
  17  * information: Portions Copyright [yyyy] [name of copyright owner]
  18  *
  19  * CDDL HEADER END
  20  */
  21 /*
  22  * Copyright 2008 Sun Microsystems, Inc.  All rights reserved.
  23  * Use is subject to license terms.
  24  */
  25 /*
  26  * Copyright 2012 Nexenta Systems, Inc. All rights reserved.
  27  */
  28 #endif /* ! codereview */

  30 /*
  31  * NL7C (Network Layer 7 Cache) as part of SOCKFS provides an in-kernel
  32  * gateway cache for the request/response message based L7 protocol HTTP
  33  * (Hypertext Transfer Protocol, see HTTP/1.1 RFC2616) in a semantically
  34  * transparent manner.
  35  *
  36  * Neither the requesting user agent (client, e.g. web browser) nor the
  37  * origin server (e.g. webserver) that provided the response cached by
  38  * NL7C are impacted in any way.
  39  *
  40  * Note, currently NL7C only processes HTTP messages via the embedded
  41  * URI of scheme http (not https nor any other), additional scheme are
  42  * intended to be supported as is practical such that much of the NL7C
  43  * framework may appear more general purpose then would be needed just
  44  * for an HTTP gateway cache.
  45  *
  46  * NL7C replaces NCA (Network Cache and Accelerator) and in the future
  47  * NCAS (NCA/SSL).
  48  *
  49  * Further, NL7C uses all NCA configuration files, see "/etc/nca/", the
  50  * NCA socket API, "AF_NCA", and "ndd /dev/nca" for backwards compatibility.
  51  */

  53 #include <sys/systm.h>
  54 #include <sys/strsun.h>
  55 #include <sys/strsubr.h>
  56 #include <inet/common.h>
  57 #include <inet/led.h>
  58 #include <inet/mi.h>
  59 #include <netinet/in.h>
  60 #include <fs/sockfs/nl7c.h>
  61 #include <fs/sockfs/nl7curi.h>
```

```
  62 #include <fs/sockfs/socktpi.h>

  64 #include <inet/nca/ncadoorhdr.h>
  65 #include <inet/nca/ncalogd.h>
  66 #include <inet/nca/ncandd.h>
  67 #include <inet/ip.h>
  68 #endif /* ! codereview */

  70 #include <sys/promif.h>

  72 /*
  73  * NL7C, NCA, NL7C logger enabled:
  74  */

  76 boolean_t        nl7c_enabled = B_FALSE;

  78 boolean_t        nl7c_logd_enabled = B_FALSE;
  79 boolean_t        nl7c_logd_started = B_FALSE;
  80 boolean_t        nl7c_logd_cycle = B_TRUE;

  82 /*
  83  * Some externs:
  84  */

  26 extern int      inet_pton(int, char *, void *);

  85 extern void     nl7c_uri_init(void);
  86 extern boolean_t nl7c_logd_init(int, caddr_t *);
  87 extern void     nl7c_nca_init(void);

  89 /*
  90  * nl7c_addr_t - a singly linked grounded list, pointed to by *nl7caddrs,
  91  * constructed at init time by parsing "/etc/nca/ncaport.conf".
  92  *
  93  * This list is searched at bind(3SOCKET) time when an application doesn't
  94  * explicitly set AF_NCA but instead uses AF_INET, if a match is found then
  95  * the underlying socket is marked sti_nl7c_flags NL7C_ENABLED.
  96  */

  98 typedef struct nl7c_addr_s {
  99         struct nl7c_addr_s *next;       /* next entry */
 100         sa_family_t     family;         /* addr type, only INET and INET6 */
 101         uint16_t        port;           /* port */
 102         union {
 103                 ipaddr_t        v4;     /* IPv4 address */
 104                 in6_addr_t      v6;     /* IPv6 address */
 105                 void            *align; /* foce alignment */
 106         }               addr;           /* address */

 108         struct sonode   *listener;      /* listen()er's sonode */
 109         boolean_t       temp;           /* temporary addr via add_addr() ? */
 110 } nl7c_addr_t;
_____unchanged_portion_omitted_

 307 /*
 308  * Inet ASCII to binary.
 309  *
 310  * Note, it's assumed that *s is a valid zero byte terminated string, and
 311  * that *p is a zero initialized struct (this is important as the value of
 312  * INADDR_ANY and IN6ADDR_ANY is zero).
 313  */

 315 static int
 316 inet_atob(char *s, nl7c_addr_t *p)
 317 {
 318         if (strcmp(s, "*") == 0) {
```

```
 319                         /* INADDR_ANY */
 320                         p->family = AF_INET;
 321                         return (0);
 322                 }
 323                 if (strcmp(s, "::") == 0) {
 324                         /* IN6ADDR_ANY */
 325                         p->family = AF_INET6;
 326                         return (0);
 327                 }
 328                 /* IPv4 address ? */
 329                 if (inet_pton(AF_INET, s, &p->addr.v4) != 1) {
 330                         /* Nop, IPv6 address ? */
 331                         if (inet_pton(AF_INET6, s, &p->addr.v6) != 1) {
 332                                 /* Nop, return error */
 333                                 return (1);
 334                         }
 335                         p->family = AF_INET6;
 336                 } else {
 337                         p->family = AF_INET;
 281                         p->addr.v4 = ntohl(p->addr.v4);
 338                 }

 340 #endif /* ! codereview */
 341         return (0);
 342 }

 344 /*
 345  * Open and read each line from "/etc/nca/ncaport.conf", the syntax of a
 346  * ncaport.conf file line is:
 347  *
 348  *      ncaport=IPaddr/Port[/Proxy]
 349  *
 350  * Where:
 351  *
 352  * ncaport - the only token recognized.
 353  *
 354  *  IPaddr - an IPv4 numeric dot address (e.g. 192.168.84.71) or '*' for
 355  *           INADDR_ANY, or an IPv6 numeric address or "::" for IN6ADDR_ANY.
 356  *
 357  *       / - IPaddr/Port separator.
 358  *
 359  *    Port - a TCP decimal port number.
 360  *
 361  * Note, all other lines will be ignored.
 362  */

 364 static void
 365 ncaportconf_read(void)
 366 {
 367         int     ret;
 368         struct vnode *vp;
 369         char    c;
 370         ssize_t resid;
 371         char    buf[1024];
 372         char    *ebp = &buf[sizeof (buf)];
 373         char    *bp = ebp;
 374         offset_t off = 0;
 375         enum parse_e {START, TOK, ADDR, PORT, EOL} parse = START;
 376         nl7c_addr_t *addrp = NULL;
 377         char    *ncaport = "ncaport";
 378         char    string[] = "XXXX:XXXX:XXXX:XXXX:XXXX:XXXX:XXXX:XXXX";
 379         char    *stringp;
 380         char    *tok;
 381         char    *portconf = "/etc/nca/ncaport.conf";

 383         ret = vn_open(portconf, UIO_SYSSPACE, FREAD, 0, &vp, 0, 0);
```

```
 384         if (ret == ENOENT) {
 385                 /* No portconf file, nothing to do */
 386                 return;
 387         }
 388         if (ret != 0) {
 389                 /* Error of some sort, tell'm about it */
 390                 cmn_err(CE_WARN, "%s: open error %d", portconf, ret);
 391                 return;
 392         }
 393         /*
 394          * Read portconf one buf[] at a time, parse one char at a time.
 395          */
 396         for (;;) {
 397                 if (bp == ebp) {
 398                         /* Nothing left in buf[], read another */
 399                         ret = vn_rdwr(UIO_READ, vp, buf, sizeof (buf), off,
 400                             UIO_SYSSPACE, 0, (rlim64_t)0, CRED(), &resid);
 401                         if (ret != 0) {
 402                                 /* Error of some sort, tell'm about it */
 403                                 cmn_err(CE_WARN, "%s: read error %d",
 404                                     portconf, ret);
 405                                 break;
 406                         }
 407                         if (resid == sizeof (buf)) {
 408                                 /* EOF, done */
 409                                 break;
 410                         }
 411                         /* Initilize per buf[] state */
 412                         bp = buf;
 413                         ebp = &buf[sizeof (buf) - resid];
 414                         off += sizeof (buf) - resid;
 415                 }
 416                 c = *bp++;
 417                 switch (parse) {
 418                 case START:
 419                         /* Initilize all per file line state */
 420                         if (addrp == NULL) {
 421                                 addrp = kmem_zalloc(sizeof (*addrp),
 422                                     KM_NOSLEEP);
 423                         }
 424                         tok = ncaport;
 425                         stringp = string;
 426                         parse = TOK;
 427                         /*FALLTHROUGH*/
 428                 case TOK:
 429                         if (c == '#') {
 430                                 /* Comment through end of line */
 431                                 parse = EOL;
 432                                 break;
 433                         }
 434                         if (isalpha(c)) {
 435                                 if (c != *tok++) {
 436                                         /* Only know one token, skip */
 437                                         parse = EOL;
 438                                 }
 439                         } else if (c == '=') {
 440                                 if (*tok != NULL) {
 441                                         /* Only know one token, skip */
 442                                         parse = EOL;
 443                                         break;
 444                                 }
 445                                 parse = ADDR;
 446                         } else if (c == '\n') {
 447                                 /* Found EOL, empty line, next line */
 448                                 parse = START;
 449                         } else {
```

```
450                             /* Unexpected char, skip */
451                             parse = EOL;
452                     }
453                     break;

455             case ADDR:
456                     if (c == '/') {
457                             /* addr/port separator, end of addr */
458                             *stringp = NULL;
459                             if (inet_atob(string, addrp)) {
460                                     /* Bad addr, skip */
461                                     parse = EOL;
462                             } else {
463                                     stringp = string;
464                                     parse = PORT;
465                             }
466                     } else {
467                             /* Save char to string */
468                             if (stringp ==
469                                 &string[sizeof (string) - 1]) {
470                                     /* Would overflow, skip */
471                                     parse = EOL;
472                             } else {
473                                     /* Copy IP addr char */
474                                     *stringp++ = c;
475                             }
476                     }
477                     break;

479             case PORT:
480                     if (isdigit(c)) {
481                             /* Save char to string */
482                             if (stringp ==
483                                 &string[sizeof (string) - 1]) {
484                                     /* Would overflow, skip */
485                                     parse = EOL;
486                             } else {
487                                     /* Copy port digit char */
488                                     *stringp++ = c;
489                             }
490                             break;
491                     } else if (c == '#' || isspace(c)) {
492                             /* End of port number, convert */
493                             *stringp = NULL;
494                             addrp->port = ntohs(atou(string));

496                             /* End of parse, add entry */
497                             nl7c_addr_add(addrp);
498                             addrp = NULL;
499                             parse = EOL;
500                     } else {
501                             /* Unrecognized char, skip */
502                             parse = EOL;
503                             break;
504                     }
505                     if (c == '\n') {
506                             /* Found EOL, start on next line */
507                             parse = START;
508                     }
509                     break;

511             case EOL:
512                     if (c == '\n') {
513                             /* Found EOL, start on next line */
514                             parse = START;
515                     }
```

```
516                     break;
517             }

519         }
520         if (addrp != NULL) {
521                 kmem_free(addrp, sizeof (*addrp));
522         }
523         (void) VOP_CLOSE(vp, FREAD, 1, (offset_t)0, CRED(), NULL);
524         VN_RELE(vp);
525 }

527 /*
528  * Open and read each line from "/etc/nca/ncakmod.conf" and parse looking
529  * for the NCA enabled, the syntax is: status=enabled, all other lines will
530  * be ignored.
531  */

533 static void
534 ncakmodconf_read(void)
535 {
536         int     ret;
537         struct vnode *vp;
538         char    c;
539         ssize_t resid;
540         char    buf[1024];
541         char    *ebp = &buf[sizeof (buf)];
542         char    *bp = ebp;
543         offset_t off = 0;
544         enum parse_e {START, TOK, EOL} parse = START;
545         char    *status = "status=enabled";
546         char    *tok;
547         char    *ncakmod = "/etc/nca/ncakmod.conf";

549         ret = vn_open(ncakmod, UIO_SYSSPACE, FREAD, 0, &vp, 0, 0);
550         if (ret == ENOENT) {
551                 /* No ncakmod file, nothing to do */
552                 return;
553         }
554         if (ret != 0) {
555                 /* Error of some sort, tell'm about it */
556                 cmn_err(CE_WARN, "%s: open error %d", status, ret);
557                 return;
558         }
559         /*
560          * Read ncakmod one buf[] at a time, parse one char at a time.
561          */
562         for (;;) {
563                 if (bp == ebp) {
564                         /* Nothing left in buf[], read another */
565                         ret = vn_rdwr(UIO_READ, vp, buf, sizeof (buf), off,
566                             UIO_SYSSPACE, 0, (rlim64_t)0, CRED(), &resid);
567                         if (ret != 0) {
568                                 /* Error of some sort, tell'm about it */
569                                 cmn_err(CE_WARN, "%s: read error %d",
570                                     status, ret);
571                                 break;
572                         }
573                         if (resid == sizeof (buf)) {
574                                 /* EOF, done */
575                                 break;
576                         }
577                         /* Initilize per buf[] state */
578                         bp = buf;
579                         ebp = &buf[sizeof (buf) - resid];
580                         off += sizeof (buf) - resid;
581                 }
```

```
582                    c = *bp++;
583                    switch (parse) {
584                    case START:
585                            /* Initilize all per file line state */
586                            tok = status;
587                            parse = TOK;
588                            /*FALLTHROUGH*/
589                    case TOK:
590                            if (c == '#') {
591                                    /* Comment through end of line */
592                                    parse = EOL;
593                                    break;
594                            }
595                            if (isalpha(c) || c == '=') {
596                                    if (c != *tok++) {
597                                            /* Only know one token, skip */
598                                            parse = EOL;
599                                    }
600                            } else if (c == '\n') {
601                                    /*
602                                     * Found EOL, if tok found done,
603                                     * else start on next-line.
604                                     */
605                                    if (*tok == NULL) {
606                                            nl7c_enabled = B_TRUE;
607                                            goto done;
608                                    }
609                                    parse = START;
610                            } else {
611                                    /* Unexpected char, skip */
612                                    parse = EOL;
613                            }
614                            break;

616                    case EOL:
617                            if (c == '\n') {
618                                    /* Found EOL, start on next line */
619                                    parse = START;
620                            }
621                            break;
622                    }
624            }
625    done:
626            (void) VOP_CLOSE(vp, FREAD, 1, (offset_t)0, CRED(), NULL);
627            VN_RELE(vp);
628    }

630    /*
631     * Open and read each line from "/etc/nca/ncalogd.conf" and parse for
632     * the tokens and token text (i.e. key and value ncalogd.conf(4)):
633     *
634     *      status=enabled
635     *
636     *      logd_file_size=[0-9]+
637     *
638     *      logd_file_name=["]filename( filename)*["]
639     */

641    static int      file_size = 1000000;
642    static caddr_t  fnv[NCA_FIOV_SZ];

644    static void
645    ncalogdconf_read(void)
646    {
647            int     ret;
```

```
648            struct vnode *vp;
649            char    c;
650            int     sz;
651            ssize_t resid;
652            char    buf[1024];
653            char    *ebp = &buf[sizeof (buf)];
654            char    *bp = ebp;
655            offset_t off = 0;
656            enum parse_e {START, TOK, TEXT, EOL} parse = START;
657            char    *tokstatus = "status\0enabled";
658            char    *toksize = "logd_file_size";
659            char    *tokfile = "logd_path_name";
660            char    *tokstatusp;
661            char    *toksizep;
662            char    *tokfilep;
663            char    *tok;
664            int     tokdelim = 0;
665            char    *ncalogd = "/etc/nca/ncalogd.conf";
666            char    *ncadeflog = "/var/nca/log";
667            char    file[TYPICALMAXPATHLEN] = {0};
668            char    *fp = file;
669            caddr_t *fnvp = fnv;

671            ret = vn_open(ncalogd, UIO_SYSSPACE, FREAD, 0, &vp, 0, 0);
672            if (ret == ENOENT) {
673                    /* No ncalogd file, nothing to do */
674                    return;
675            }
676            if (ret != 0) {
677                    /* Error of some sort, tell'm about it */
678                    cmn_err(CE_WARN, "ncalogdconf_read: %s: open error(%d).",
679                        ncalogd, ret);
680                    return;
681            }
682            /*
683             * Read ncalogd.conf one buf[] at a time, parse one char at a time.
684             */
685            for (;;) {
686                    if (bp == ebp) {
687                            /* Nothing left in buf[], read another */
688                            ret = vn_rdwr(UIO_READ, vp, buf, sizeof (buf), off,
689                                UIO_SYSSPACE, 0, (rlim64_t)0, CRED(), &resid);
690                            if (ret != 0) {
691                                    /* Error of some sort, tell'm about it */
692                                    cmn_err(CE_WARN, "%s: read error %d",
693                                        ncalogd, ret);
694                                    break;
695                            }
696                            if (resid == sizeof (buf)) {
697                                    /* EOF, done */
698                                    break;
699                            }
700                            /* Initilize per buf[] state */
701                            bp = buf;
702                            ebp = &buf[sizeof (buf) - resid];
703                            off += sizeof (buf) - resid;
704                    }
705                    c = *bp++;
706                    switch (parse) {
707                    case START:
708                            /* Initilize all per file line state */
709                            tokstatusp = tokstatus;
710                            toksizep = toksize;
711                            tokfilep = tokfile;
712                            tok = NULL;
713                            parse = TOK;
```

```
 714                                         sz = 0;
 715                                 /*FALLTHROUGH*/
 716                         case TOK:
 717                                 if (isalpha(c) || c == '_') {
 718                                         /*
 719                                          * Found a valid tok char, if matches
 720                                          * any of the tokens continue else NULL
 721                                          * then string pointer.
 722                                          */
 723                                         if (tokstatusp != NULL && c != *tokstatusp++)
 724                                                 tokstatusp = NULL;
 725                                         if (toksizep != NULL && c != *toksizep++)
 726                                                 toksizep = NULL;
 727                                         if (tokfilep != NULL && c != *tokfilep++)
 728                                                 tokfilep = NULL;

 730                                         if (tokstatusp == NULL &&
 731                                             toksizep == NULL &&
 732                                             tokfilep == NULL) {
 733                                                 /*
 734                                                  * All tok string pointers are NULL
 735                                                  * so skip rest of line.
 736                                                  */
 737                                                 parse = EOL;
 738                                         }
 739                                 } else if (c == '=') {
 740                                         /*
 741                                          * Found tok separator, if tok found get
 742                                          * tok text, else skip rest of line.
 743                                          */
 744                                         if (tokstatusp != NULL && *tokstatusp == NULL)
 745                                                 tok = tokstatus;
 746                                         else if (toksizep != NULL && *toksizep == NULL)
 747                                                 tok = toksize;
 748                                         else if (tokfilep != NULL && *tokfilep == NULL)
 749                                                 tok = tokfile;
 750                                         if (tok != NULL)
 751                                                 parse = TEXT;
 752                                         else
 753                                                 parse = EOL;
 754                                 } else if (c == '\n') {
 755                                         /* Found EOL, start on next line */
 756                                         parse = START;
 757                                 } else {
 758                                         /* Comment or unknown char, skip rest of line */
 759                                         parse = EOL;
 760                                 }
 761                                 break;
 762                         case TEXT:
 763                                 if (c == '\n') {
 764                                         /*
 765                                          * Found EOL, finish up tok text processing
 766                                          * (if any) and start on next line.
 767                                          */
 768                                         if (tok == tokstatus) {
 769                                                 if (*++tokstatusp == NULL)
 770                                                         nl7c_logd_enabled = B_TRUE;
 771                                         } else if (tok == toksize) {
 772                                                 file_size = sz;
 773                                         } else if (tok == tokfile) {
 774                                                 if (tokdelim == 0) {
 775                                                         /* Non delimited path name */
 776                                                         *fnvp++ = strdup(file);
 777                                                 } else if (fp != file) {
 778                                                         /* No closing delimiter */
 779                                                         /*EMPTY*/;
```

```
 780                                                 }
 781                                         }
 782                                         parse = START;
 783                                 } else if (tok == tokstatus) {
 784                                         if (! isalpha(c) || *++tokstatusp == NULL ||
 785                                             c != *tokstatusp) {
 786                                                 /* Not enabled, skip line */
 787                                                 parse = EOL;
 788                                         }
 789                                 } else if (tok == toksize) {
 790                                         if (isdigit(c)) {
 791                                                 sz *= 10;
 792                                                 sz += c - '0';
 793                                         } else {
 794                                                 /* Not a decimal digit, skip line */
 795                                                 parse = EOL;
 796                                         }
 797                                 } else {
 798                                         /* File name */
 799                                         if (c == '"' && tokdelim++ == 0) {
 800                                                 /* Opening delimiter, skip */
 801                                                 /*EMPTY*/;
 802                                         } else if (c == '"' || c == ' ') {
 803                                                 /* List delim or filename separator */
 804                                                 *fnvp++ = strdup(file);
 805                                                 fp = file;
 806                                         } else if (fp < &file[sizeof (file) - 1]) {
 807                                                 /* Filename char */
 808                                                 *fp++ = c;
 809                                         } else {
 810                                                 /* Filename to long, skip line */
 811                                                 parse = EOL;
 812                                         }
 813                                 }
 814                                 break;

 816                         case EOL:
 817                                 if (c == '\n') {
 818                                         /* Found EOL, start on next line */
 819                                         parse = START;
 820                                 }
 821                                 break;
 822                         }

 824         }
 825 done:
 826         (void) VOP_CLOSE(vp, FREAD, 1, (offset_t)0, CRED(), NULL);
 827         VN_RELE(vp);

 829         if (nl7c_logd_enabled) {
 830                 if (fnvp == fnv) {
 831                         /*
 832                          * No logfile was specified and found so
 833                          * so use defualt NCA log file path.
 834                          */
 835                         *fnvp++ = strdup(ncadeflog);
 836                 }
 837                 if (fnvp < &fnv[NCA_FIOV_SZ]) {
 838                         /* NULL terminate list */
 839                         *fnvp = NULL;
 840                 }
 841         }
 842 }

 844 void
 845 nl7clogd_startup(void)
```

```
846 {
847         static kmutex_t startup;

849         /*
850          * Called on the first log() attempt, have to wait until then to
851          * initialize logd as at logdconf_read() the root fs is read-only.
852          */
853         mutex_enter(&startup);
854         if (nl7c_logd_started) {
855                 /* Lost the race, nothing todo */
856                 mutex_exit(&startup);
857                 return;
858         }
859         nl7c_logd_started = B_TRUE;
860         if (! nl7c_logd_init(file_size, fnv)) {
861                 /* Failure, disable logging */
862                 nl7c_logd_enabled = B_FALSE;
863                 cmn_err(CE_WARN, "nl7clogd_startup: failed, disabling loggin");
864                 mutex_exit(&startup);
865                 return;
866         }
867         mutex_exit(&startup);
868 }


871 void
872 nl7c_startup()
873 {
874         /*
875          * Open, read, and parse the NCA logd configuration file,
876          * then initialize URI processing and NCA compat.
877          */
878         ncalogdconf_read();
879         nl7c_uri_init();
880         nl7c_nca_init();
881 }

883 void
884 nl7c_init()
885 {
886         /* Open, read, and parse the NCA kmod configuration file */
887         ncakmodconf_read();

889         if (nl7c_enabled) {
890                 /*
891                  * NL7C is enabled so open, read, and parse
892                  * the NCA address/port configuration file
893                  * and call startup() to finish config/init.
894                  */
895                 ncaportconf_read();
896                 nl7c_startup();
897         }
898 }

900 /*
901  * The main processing function called by accept() on a newly created
902  * socket prior to returning it to the caller of accept().
903  *
904  * Here data is read from the socket until a completed L7 request parse
905  * is completed. Data will be read in the context of the user thread
906  * which called accept(), when parse has been completed either B_TRUE
907  * or B_FALSE will be returned.
908  *
909  * If NL7C successfully process the L7 protocol request, i.e. generates
910  * a response, B_TRUE will be returned.
911  *
```

```
912  * Else, B_FALSE will be returned if NL7C can't process the request:
913  *
914  * 1) Couldn't locate a URI within the request.
915  *
916  * 2) URI scheme not reqcognized.
917  *
918  * 3) A request which can't be processed.
919  *
920  * 4) A request which could be processed but NL7C dosen't currently have
921  *    the response data. In which case NL7C will parse the returned response
922  *    from the application for possible caching for subsequent request(s).
923  */

925 volatile uint64_t nl7c_proc_cnt = 0;
926 volatile uint64_t nl7c_proc_error = 0;
927 volatile uint64_t nl7c_proc_ETIME = 0;
928 volatile uint64_t nl7c_proc_again = 0;
929 volatile uint64_t nl7c_proc_next = 0;
930 volatile uint64_t nl7c_proc_rcv = 0;
931 volatile uint64_t nl7c_proc_noLRI = 0;
932 volatile uint64_t nl7c_proc_nodata = 0;
933 volatile uint64_t nl7c_proc_parse = 0;

935 boolean_t
936 nl7c_process(struct sonode *so, boolean_t nonblocking)
937 {
938         vnode_t *vp = SOTOV(so);
939         sotpi_info_t *sti = SOTOTPI(so);
940         mblk_t  *rmp = sti->sti_nl7c_rcv_mp;
941         clock_t timout;
942         rval_t  rval;
943         uchar_t pri;
944         int     pflag;
945         int     error;
946         boolean_t more;
947         boolean_t ret = B_FALSE;
948         boolean_t first = B_TRUE;
949         boolean_t pollin = (sti->sti_nl7c_flags & NL7C_POLLIN);

951         nl7c_proc_cnt++;

953         /* Caller has so_lock enter()ed */
954         error = so_lock_read_intr(so, nonblocking ? FNDELAY|FNONBLOCK : 0);
955         if (error) {
956                 /* Couldn't read lock, pass on this socket */
957                 sti->sti_nl7c_flags = 0;
958                 nl7c_proc_noLRI++;
959                 return (B_FALSE);
960         }
961         /* Exit so_lock for now, will be reenter()ed prior to return */
962         mutex_exit(&so->so_lock);

964         if (pollin)
965                 sti->sti_nl7c_flags &= ~NL7C_POLLIN;

967         /* Initialize some kstrgetmsg() constants */
968         pflag = MSG_ANY | MSG_DELAYERROR;
969         pri = 0;
970         if (nonblocking) {
971                 /* Non blocking so don't block */
972                 timout = 0;
973         } else if (sti->sti_nl7c_flags & NL7C_SOPERSIST) {
974                 /* 2nd or more time(s) here so use keep-alive value */
975                 timout = nca_http_keep_alive_timeout;
976         } else {
977                 /* 1st time here so use connection value */
```

```
 978                             timout = nca_http_timeout;
 979                     }

 981             rval.r_vals = 0;
 982             do {
 983                             /*
 984                              * First time through, if no data left over from a previous
 985                              * kstrgetmsg() then try to get some, else just process it.
 986                              *
 987                              * Thereafter, rmp = NULL after the successful kstrgetmsg()
 988                              * so try to get some new data and append to list (i.e. until
 989                              * enough fragments are collected for a successful parse).
 990                              */
 991                             if (rmp == NULL) {

 993                                     error = kstrgetmsg(vp, &rmp, NULL, &pri, &pflag,
 994                                         timout, &rval);
 995                                     if (error) {
 996                                             if (error == ETIME) {
 997                                                     /* Timeout */
 998                                                     nl7c_proc_ETIME++;
 999                                             } else if (error != EWOULDBLOCK) {
1000                                                     /* Error of some sort */
1001                                                     nl7c_proc_error++;
1002                                                     rval.r_v.r_v2 = error;
1003                                                     sti->sti_nl7c_flags = 0;
1004                                                     break;
1005                                             }
1006                                             error = 0;
1007                                     }
1008                                     if (rmp != NULL) {
1009                                             mblk_t  *mp = sti->sti_nl7c_rcv_mp;

1012                                             if (mp == NULL) {
1013                                                     /* Just new data, common case */
1014                                                     sti->sti_nl7c_rcv_mp = rmp;
1015                                             } else {
1016                                                     /* Add new data to tail */
1017                                                     while (mp->b_cont != NULL)
1018                                                             mp = mp->b_cont;
1019                                                     mp->b_cont = rmp;
1020                                             }
1021                                     }
1022                                     if (sti->sti_nl7c_rcv_mp == NULL) {
1023                                             /* No data */
1024                                             nl7c_proc_nodata++;
1025                                             if (timout > 0 || (first && pollin)) {
1026                                                     /* Expected data so EOF */
1027                                                     ret = B_TRUE;
1028                                             } else if (sti->sti_nl7c_flags &
1029                                                 NL7C_SOPERSIST) {
1030                                                     /* Persistent so just checking */
1031                                                     ret = B_FALSE;
1032                                             }
1033                                             break;
1034                                     }
1035                                     rmp = NULL;
1036                             }
1037                     first = B_FALSE;
1038             again:
1039                     nl7c_proc_parse++;

1041                     more = nl7c_parse(so, nonblocking, &ret);

1043                     if (ret == B_TRUE && (sti->sti_nl7c_flags & NL7C_SOPERSIST)) {
```

```
1044                             /*
1045                              * Parse complete, cache hit, response on its way,
1046                              * socket is persistent so try to process the next
1047                              * request.
1048                              */
1049                             if (nonblocking) {
1050                                     ret = B_FALSE;
1051                                     break;
1052                             }
1053                             if (sti->sti_nl7c_rcv_mp) {
1054                                     /* More recv-side data, pipelined */
1055                                     nl7c_proc_again++;
1056                                     goto again;
1057                             }
1058                             nl7c_proc_next++;
1059                             if (nonblocking)
1060                                     timout = 0;
1061                             else
1062                                     timout = nca_http_keep_alive_timeout;

1064                             more = B_TRUE;
1065                     }

1067             } while (more);

1069             if (sti->sti_nl7c_rcv_mp) {
1070                     nl7c_proc_rcv++;
1071             }
1072             sti->sti_nl7c_rcv_rval = rval.r_vals;
1073             /* Renter so_lock, caller called with it enter()ed */
1074             mutex_enter(&so->so_lock);
1075             so_unlock_read(so);

1077             return (ret);
1078 }
```

```
   1 /*
   2  * CDDL HEADER START
   3  *
   4  * The contents of this file are subject to the terms of the
   5  * Common Development and Distribution License (the "License").
   6  * You may not use this file except in compliance with the License.
   7  *
   8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
   9  * or http://www.opensolaris.org/os/licensing.
  10  * See the License for the specific language governing permissions
  11  * and limitations under the License.
  12  *
  13  * When distributing Covered Code, include this CDDL HEADER in each
  14  * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
  15  * If applicable, add the following below this CDDL HEADER, with the
  16  * fields enclosed by brackets "[]" replaced with your own identifying
  17  * information: Portions Copyright [yyyy] [name of copyright owner]
  18  *
  19  * CDDL HEADER END
  20  */

  22 /*
  23  * Copyright (c) 1991, 2010, Oracle and/or its affiliates. All rights reserved.
  24  * Copyright (c) 1990 Mentat Inc.
  25  */
  26 /*
  27  * Copyright 2012 Nexenta Systems, Inc. All rights reserved.
  28  */
  29 #endif /* ! codereview */

  31 #ifndef _INET_IP_H
  32 #define _INET_IP_H

  34 #ifdef  __cplusplus
  35 extern "C" {
  36 #endif

  38 #include <sys/isa_defs.h>
  39 #include <sys/types.h>
  40 #include <inet/mib2.h>
  41 #include <inet/nd.h>
  42 #include <sys/atomic.h>
  43 #include <net/if_dl.h>
  44 #include <net/if.h>
  45 #include <netinet/ip.h>
  46 #include <netinet/igmp.h>
  47 #include <sys/neti.h>
  48 #include <sys/hook.h>
  49 #include <sys/hook_event.h>
  50 #include <sys/hook_impl.h>
  51 #include <inet/ip_stack.h>

  53 #ifdef _KERNEL
  54 #include <netinet/ip6.h>
  55 #include <sys/avl.h>
  56 #include <sys/list.h>
  57 #include <sys/vmem.h>
  58 #include <sys/squeue.h>
  59 #include <sys/systm.h>
  60 #include <net/route.h>
  61 #include <net/radix.h>
```

```
  62 #include <sys/modhash.h>

  64 #ifdef DEBUG
  65 #define CONN_DEBUG
  66 #endif

  68 #define IP_DEBUG
  69 /*
  70  * The mt-streams(9F) flags for the IP module; put here so that other
  71  * "drivers" that are actually IP (e.g., ICMP, UDP) can use the same set
  72  * of flags.
  73  */
  74 #define IP_DEVMTFLAGS D_MP
  75 #endif  /* _KERNEL */

  77 #define IP_MOD_NAME     "ip"
  78 #define IP_DEV_NAME     "/dev/ip"
  79 #define IP6_DEV_NAME    "/dev/ip6"

  81 #define UDP_MOD_NAME    "udp"
  82 #define UDP_DEV_NAME    "/dev/udp"
  83 #define UDP6_DEV_NAME   "/dev/udp6"

  85 #define TCP_MOD_NAME    "tcp"
  86 #define TCP_DEV_NAME    "/dev/tcp"
  87 #define TCP6_DEV_NAME   "/dev/tcp6"

  89 #define SCTP_MOD_NAME   "sctp"

  91 #ifndef _IPADDR_T
  92 #define _IPADDR_T
  93 typedef uint32_t ipaddr_t;
  94 #endif

  96 /* Number of bits in an address */
  97 #define IP_ABITS                32
  98 #define IPV4_ABITS              IP_ABITS
  99 #define IPV6_ABITS              128
 100 #define IP_MAX_HW_LEN   40

 102 #define IP_HOST_MASK            (ipaddr_t)0xffffffffU

 104 #define IP_CSUM(mp, off, sum)           (~ip_cksum(mp, off, sum) & 0xFFFF)
 105 #define IP_CSUM_PARTIAL(mp, off, sum)   ip_cksum(mp, off, sum)
 106 #define IP_BCSUM_PARTIAL(bp, len, sum)  bcksum(bp, len, sum)

 108 #define ILL_FRAG_HASH_TBL_COUNT ((unsigned int)64)
 109 #define ILL_FRAG_HASH_TBL_SIZE  (ILL_FRAG_HASH_TBL_COUNT * sizeof (ipfb_t))

 111 #define IPV4_ADDR_LEN                   4
 112 #define IP_ADDR_LEN                     IPV4_ADDR_LEN
 113 #define IP_ARP_PROTO_TYPE               0x0800

 115 #define IPV4_VERSION                    4
 116 #define IP_VERSION                      IPV4_VERSION
 117 #define IP_SIMPLE_HDR_LENGTH_IN_WORDS   5
 118 #define IP_SIMPLE_HDR_LENGTH            20
 119 #define IP_MAX_HDR_LENGTH               60

 121 #define IP_MAX_OPT_LENGTH (IP_MAX_HDR_LENGTH-IP_SIMPLE_HDR_LENGTH)

 123 #define IP_MIN_MTU                      (IP_MAX_HDR_LENGTH + 8) /* 68 bytes */

 125 /*
 126  * XXX IP_MAXPACKET is defined in <netinet/ip.h> as well. At some point the
 127  * 2 files should be cleaned up to remove all redundant definitions.
```

```
 128 */
 129 #define IP_MAXPACKET                    65535
 130 #define IP_SIMPLE_HDR_VERSION \
 131         ((IP_VERSION << 4) | IP_SIMPLE_HDR_LENGTH_IN_WORDS)

 133 #define UDPH_SIZE                       8

 135 /*
 136  * Constants and type definitions to support IP IOCTL commands
 137  */
 138 #define IP_IOCTL                        (('i'<<8)|'p')
 139 #define IP_IOC_IRE_DELETE               4
 140 #define IP_IOC_IRE_DELETE_NO_REPLY      5
 141 #define IP_IOC_RTS_REQUEST              7

 143 /* Common definitions used by IP IOCTL data structures */
 144 typedef struct ipllcmd_s {
 145         uint_t  ipllc_cmd;
 146         uint_t  ipllc_name_offset;
 147         uint_t  ipllc_name_length;
 148 } ipllc_t;

 150 /* IP IRE Delete Command Structure. */
 151 typedef struct ipid_s {
 152         ipllc_t ipid_ipllc;
 153         uint_t  ipid_ire_type;
 154         uint_t  ipid_addr_offset;
 155         uint_t  ipid_addr_length;
 156         uint_t  ipid_mask_offset;
 157         uint_t  ipid_mask_length;
 158 } ipid_t;

 160 #define ipid_cmd                ipid_ipllc.ipllc_cmd

 162 #ifdef _KERNEL
 163 /*
 164  * Temporary state for ip options parser.
 165  */
 166 typedef struct ipoptp_s
 167 {
 168         uint8_t         *ipoptp_next;   /* next option to look at */
 169         uint8_t         *ipoptp_end;    /* end of options */
 170         uint8_t         *ipoptp_cur;    /* start of current option */
 171         uint8_t         ipoptp_len;     /* length of current option */
 172         uint32_t        ipoptp_flags;
 173 } ipoptp_t;

 175 /*
 176  * Flag(s) for ipoptp_flags
 177  */
 178 #define IPOPTP_ERROR    0x00000001
 179 #endif  /* _KERNEL */

 181 /* Controls forwarding of IP packets, set via ipadm(1M)/ndd(1M) */
 182 #define IP_FORWARD_NEVER        0
 183 #define IP_FORWARD_ALWAYS       1

 185 #define WE_ARE_FORWARDING(ipst) ((ipst)->ips_ip_forwarding == IP_FORWARD_ALWAYS)

 187 #define IPH_HDR_LENGTH(ipha)                                            \
 188         ((int)(((ipha_t *)ipha)->ipha_version_and_hdr_length & 0xF) << 2)

 190 #define IPH_HDR_VERSION(ipha)                                           \
 191         ((int)(((ipha_t *)ipha)->ipha_version_and_hdr_length) >> 4)

 193 #ifdef _KERNEL
```

```
 194 /*
 195  * IP reassembly macros.  We hide starting and ending offsets in b_next and
 196  * b_prev of messages on the reassembly queue.  The messages are chained using
 197  * b_cont.  These macros are used in ip_reassemble() so we don't have to see
 198  * the ugly casts and assignments.
 199  * Note that the offsets are <= 64k i.e. a uint_t is sufficient to represent
 200  * them.
 201  */
 202 #define IP_REASS_START(mp)              ((uint_t)(uintptr_t)((mp)->b_next))
 203 #define IP_REASS_SET_START(mp, u)       \
 204         ((mp)->b_next = (mblk_t *)(uintptr_t)(u))
 205 #define IP_REASS_END(mp)                ((uint_t)(uintptr_t)((mp)->b_prev))
 206 #define IP_REASS_SET_END(mp, u)         \
 207         ((mp)->b_prev = (mblk_t *)(uintptr_t)(u))

 209 #define IP_REASS_COMPLETE       0x1
 210 #define IP_REASS_PARTIAL        0x2
 211 #define IP_REASS_FAILED         0x4

 213 /*
 214  * Test to determine whether this is a module instance of IP or a
 215  * driver instance of IP.
 216  */
 217 #define CONN_Q(q)       (WR(q)->q_next == NULL)

 219 #define Q_TO_CONN(q)    ((conn_t *)(q)->q_ptr)
 220 #define Q_TO_TCP(q)     (Q_TO_CONN((q))->conn_tcp)
 221 #define Q_TO_UDP(q)     (Q_TO_CONN((q))->conn_udp)
 222 #define Q_TO_ICMP(q)    (Q_TO_CONN((q))->conn_icmp)
 223 #define Q_TO_RTS(q)     (Q_TO_CONN((q))->conn_rts)

 225 #define CONNP_TO_WQ(connp)      ((connp)->conn_wq)
 226 #define CONNP_TO_RQ(connp)      ((connp)->conn_rq)

 228 #define GRAB_CONN_LOCK(q)       {                                       \
 229         if (q != NULL && CONN_Q(q))                                     \
 230                 mutex_enter(&(Q_TO_CONN(q))->conn_lock);                \
 231 }

 233 #define RELEASE_CONN_LOCK(q)    {                                       \
 234         if (q != NULL && CONN_Q(q))                                     \
 235                 mutex_exit(&(Q_TO_CONN(q))->conn_lock);                 \
 236 }

 238 /*
 239  * Ref counter macros for ioctls. This provides a guard for TCP to stop
 240  * tcp_close from removing the rq/wq whilst an ioctl is still in flight on the
 241  * stream. The ioctl could have been queued on e.g. an ipsq. tcp_close will wait
 242  * until the ioctlref count is zero before proceeding.
 243  * Ideally conn_oper_pending_ill would be used for this purpose. However, in the
 244  * case where an ioctl is aborted or interrupted, it can be cleared prematurely.
 245  * There are also some race possibilities between ip and the stream head which
 246  * can also end up with conn_oper_pending_ill being cleared prematurely. So, to
 247  * avoid these situations, we use a dedicated ref counter for ioctls which is
 248  * used in addition to and in parallel with the normal conn_ref count.
 249  */
 250 #define CONN_INC_IOCTLREF_LOCKED(connp) {                               \
 251         ASSERT(MUTEX_HELD(&(connp)->conn_lock));                        \
 252         DTRACE_PROBE1(conn__inc__ioctlref, conn_t *, (connp));          \
 253         (connp)->conn_ioctlref++;                                       \
 254         mutex_exit(&(connp)->conn_lock);                                \
 255 }

 257 #define CONN_INC_IOCTLREF(connp)        {                               \
 258         mutex_enter(&(connp)->conn_lock);                               \
 259         CONN_INC_IOCTLREF_LOCKED(connp);                                \
```

```
 260 }

 262 #define CONN_DEC_IOCTLREF(connp)             {                          \
 263         mutex_enter(&(connp)->conn_lock);                              \
 264         DTRACE_PROBE1(conn__dec__ioctlref, conn_t *, (connp));  \
 265         /* Make sure conn_ioctlref will not underflow. */              \
 266         ASSERT((connp)->conn_ioctlref != 0);                           \
 267         if ((--(connp)->conn_ioctlref == 0) &&                         \
 268             ((connp)->conn_state_flags & CONN_CLOSING)) {              \
 269                 cv_broadcast(&(connp)->conn_cv);                       \
 270         }                                                              \
 271         mutex_exit(&(connp)->conn_lock);                               \
 272 }


 275 /*
 276  * Complete the pending operation. Usually an ioctl. Can also
 277  * be a bind or option management request that got enqueued
 278  * in an ipsq_t. Called on completion of the operation.
 279  */
 280 #define CONN_OPER_PENDING_DONE(connp)     {                             \
 281         mutex_enter(&(connp)->conn_lock);                              \
 282         (connp)->conn_oper_pending_ill = NULL;                         \
 283         cv_broadcast(&(connp)->conn_refcv);                            \
 284         mutex_exit(&(connp)->conn_lock);                               \
 285         CONN_DEC_REF(connp);                                           \
 286 }

 288 /*
 289  * Values for squeue switch:
 290  */
 291 #define IP_SQUEUE_ENTER_NODRAIN 1
 292 #define IP_SQUEUE_ENTER 2
 293 #define IP_SQUEUE_FILL 3

 295 extern int ip_squeue_flag;

 297 /* IP Fragmentation Reassembly Header */
 298 typedef struct ipf_s {
 299         struct ipf_s    *ipf_hash_next;
 300         struct ipf_s    **ipf_ptphn;    /* Pointer to previous hash next. */
 301         uint32_t        ipf_ident;      /* Ident to match. */
 302         uint8_t         ipf_protocol;   /* Protocol to match. */
 303         uchar_t         ipf_last_frag_seen : 1; /* Last fragment seen ? */
 304         time_t          ipf_timestamp;  /* Reassembly start time. */
 305         mblk_t          *ipf_mp;        /* mblk we live in. */
 306         mblk_t          *ipf_tail_mp;   /* Frag queue tail pointer. */
 307         int             ipf_hole_cnt;   /* Number of holes (hard-case). */
 308         int             ipf_end;        /* Tail end offset (0 -> hard-case). */
 309         uint_t          ipf_gen;        /* Frag queue generation */
 310         size_t          ipf_count;      /* Count of bytes used by frag */
 311         uint_t          ipf_nf_hdr_len; /* Length of nonfragmented header */
 312         in6_addr_t      ipf_v6src;      /* IPv6 source address */
 313         in6_addr_t      ipf_v6dst;      /* IPv6 dest address */
 314         uint_t          ipf_prev_nexthdr_offset; /* Offset for nexthdr value */
 315         uint8_t         ipf_ecn;        /* ECN info for the fragments */
 316         uint8_t         ipf_num_dups;   /* Number of times dup frags recvd */
 317         uint16_t        ipf_checksum_flags; /* Hardware checksum flags */
 318         uint32_t        ipf_checksum;   /* Partial checksum of fragment data */
 319 } ipf_t;

 321 /*
 322  * IPv4 Fragments
 323  */
 324 #define IS_V4_FRAGMENT(ipha_fragment_offset_and_flags)                 \
 325         (((ntohs(ipha_fragment_offset_and_flags) & IPH_OFFSET) != 0) || \
```

```
 326         ((ntohs(ipha_fragment_offset_and_flags) & IPH_MF) != 0))

 328 #define ipf_src V4_PART_OF_V6(ipf_v6src)
 329 #define ipf_dst V4_PART_OF_V6(ipf_v6dst)

 331 #endif /* _KERNEL */

 333 /* ICMP types */
 334 #define ICMP_ECHO_REPLY                 0
 335 #define ICMP_DEST_UNREACHABLE           3
 336 #define ICMP_SOURCE_QUENCH              4
 337 #define ICMP_REDIRECT                   5
 338 #define ICMP_ECHO_REQUEST               8
 339 #define ICMP_ROUTER_ADVERTISEMENT       9
 340 #define ICMP_ROUTER_SOLICITATION        10
 341 #define ICMP_TIME_EXCEEDED              11
 342 #define ICMP_PARAM_PROBLEM              12
 343 #define ICMP_TIME_STAMP_REQUEST         13
 344 #define ICMP_TIME_STAMP_REPLY           14
 345 #define ICMP_INFO_REQUEST               15
 346 #define ICMP_INFO_REPLY                 16
 347 #define ICMP_ADDRESS_MASK_REQUEST       17
 348 #define ICMP_ADDRESS_MASK_REPLY         18

 350 /* Evaluates to true if the ICMP type is an ICMP error */
 351 #define ICMP_IS_ERROR(type)     (                      \
 352         (type) == ICMP_DEST_UNREACHABLE ||      \
 353         (type) == ICMP_SOURCE_QUENCH ||         \
 354         (type) == ICMP_TIME_EXCEEDED ||         \
 355         (type) == ICMP_PARAM_PROBLEM)

 357 /* ICMP_TIME_EXCEEDED codes */
 358 #define ICMP_TTL_EXCEEDED               0
 359 #define ICMP_REASSEMBLY_TIME_EXCEEDED   1

 361 /* ICMP_DEST_UNREACHABLE codes */
 362 #define ICMP_NET_UNREACHABLE            0
 363 #define ICMP_HOST_UNREACHABLE           1
 364 #define ICMP_PROTOCOL_UNREACHABLE       2
 365 #define ICMP_PORT_UNREACHABLE           3
 366 #define ICMP_FRAGMENTATION_NEEDED       4
 367 #define ICMP_SOURCE_ROUTE_FAILED        5
 368 #define ICMP_DEST_NET_UNKNOWN           6
 369 #define ICMP_DEST_HOST_UNKNOWN          7
 370 #define ICMP_SRC_HOST_ISOLATED          8
 371 #define ICMP_DEST_NET_UNREACH_ADMIN     9
 372 #define ICMP_DEST_HOST_UNREACH_ADMIN    10
 373 #define ICMP_DEST_NET_UNREACH_TOS       11
 374 #define ICMP_DEST_HOST_UNREACH_TOS      12

 376 /* ICMP Header Structure */
 377 typedef struct icmph_s {
 378         uint8_t         icmph_type;
 379         uint8_t         icmph_code;
 380         uint16_t        icmph_checksum;
 381         union {
 382                 struct { /* ECHO request/response structure */
 383                         uint16_t        u_echo_ident;
 384                         uint16_t        u_echo_seqnum;
 385                 } u_echo;
 386                 struct { /* Destination unreachable structure */
 387                         uint16_t        u_du_zero;
 388                         uint16_t        u_du_mtu;
 389                 } u_du;
 390                 struct { /* Parameter problem structure */
 391                         uint8_t         u_pp_ptr;
```

```
392                            uint8_t        u_pp_rsvd[3];
393                    } u_pp;
394                    struct { /* Redirect structure */
395                            ipaddr_t       u_rd_gateway;
396                    } u_rd;
397            } icmph_u;
398 } icmph_t;

400 #define icmph_echo_ident        icmph_u.u_echo.u_echo_ident
401 #define icmph_echo_seqnum       icmph_u.u_echo.u_echo_seqnum
402 #define icmph_du_zero           icmph_u.u_du.u_du_zero
403 #define icmph_du_mtu            icmph_u.u_du.u_du_mtu
404 #define icmph_pp_ptr            icmph_u.u_pp.u_pp_ptr
405 #define icmph_rd_gateway        icmph_u.u_rd.u_rd_gateway

407 #define ICMPH_SIZE      8

409 /*
410  * Minimum length of transport layer header included in an ICMP error
411  * message for it to be considered valid.
412  */
413 #define ICMP_MIN_TP_HDR_LEN     8

415 /* Aligned IP header */
416 typedef struct ipha_s {
417         uint8_t     ipha_version_and_hdr_length;
418         uint8_t     ipha_type_of_service;
419         uint16_t    ipha_length;
420         uint16_t    ipha_ident;
421         uint16_t    ipha_fragment_offset_and_flags;
422         uint8_t     ipha_ttl;
423         uint8_t     ipha_protocol;
424         uint16_t    ipha_hdr_checksum;
425         ipaddr_t    ipha_src;
426         ipaddr_t    ipha_dst;
427 } ipha_t;

429 /*
430  * IP Flags
431  *
432  * Some of these constant names are copied for the DTrace IP provider in
433  * usr/src/lib/libdtrace/common/{ip.d.in, ip.sed.in}, which should be kept
434  * in sync.
435  */
436 #define IPH_DF          0x4000  /* Don't fragment */
437 #define IPH_MF          0x2000  /* More fragments to come */
438 #define IPH_OFFSET      0x1FFF  /* Where the offset lives */

440 /* Byte-order specific values */
441 #ifdef  _BIG_ENDIAN
442 #define IPH_DF_HTONS     0x4000  /* Don't fragment */
443 #define IPH_MF_HTONS     0x2000  /* More fragments to come */
444 #define IPH_OFFSET_HTONS 0x1FFF /* Where the offset lives */
445 #else
446 #define IPH_DF_HTONS     0x0040  /* Don't fragment */
447 #define IPH_MF_HTONS     0x0020  /* More fragments to come */
448 #define IPH_OFFSET_HTONS 0xFF1F /* Where the offset lives */
449 #endif

451 /* ECN code points for IPv4 TOS byte and IPv6 traffic class octet. */
452 #define IPH_ECN_NECT    0x0     /* Not ECN-Capable Transport */
453 #define IPH_ECN_ECT1    0x1     /* ECN-Capable Transport, ECT(1) */
454 #define IPH_ECN_ECT0    0x2     /* ECN-Capable Transport, ECT(0) */
455 #define IPH_ECN_CE      0x3     /* ECN-Congestion Experienced (CE) */

457 struct ill_s;
```

```
459 typedef void ip_v6intfid_func_t(struct ill_s *, in6_addr_t *);
460 typedef void ip_v6mapinfo_func_t(struct ill_s *, uchar_t *, uchar_t *);
461 typedef void ip_v4mapinfo_func_t(struct ill_s *, uchar_t *, uchar_t *);

463 /* IP Mac info structure */
464 typedef struct ip_m_s {
465         t_uscalar_t             ip_m_mac_type;  /* From <sys/dlpi.h> */
466         int                     ip_m_type;      /* From <net/if_types.h> */
467         t_uscalar_t             ip_m_ipv4sap;
468         t_uscalar_t             ip_m_ipv6sap;
469         ip_v4mapinfo_func_t     *ip_m_v4mapping;
470         ip_v6mapinfo_func_t     *ip_m_v6mapping;
471         ip_v6intfid_func_t      *ip_m_v6intfid;
472         ip_v6intfid_func_t      *ip_m_v6destintfid;
473 } ip_m_t;

475 /*
476  * The following functions attempt to reduce the link layer dependency
477  * of the IP stack. The current set of link specific operations are:
478  * a. map from IPv4 class D (224.0/4) multicast address range or the
479  * IPv6 multicast address range (ff00::/8) to the link layer multicast
480  * address.
481  * b. derive the default IPv6 interface identifier from the interface.
482  * c. derive the default IPv6 destination interface identifier from
483  * the interface (point-to-point only).
484  */
485 extern  void ip_mcast_mapping(struct ill_s *, uchar_t *, uchar_t *);
486 /* ip_m_v6*intfid return void and are never NULL */
487 #define MEDIA_V6INTFID(ip_m, ill, v6ptr) (ip_m)->ip_m_v6intfid(ill, v6ptr)
488 #define MEDIA_V6DESTINTFID(ip_m, ill, v6ptr) \
489         (ip_m)->ip_m_v6destintfid(ill, v6ptr)

491 /* Router entry types */
492 #define IRE_BROADCAST           0x0001  /* Route entry for broadcast address */
493 #define IRE_DEFAULT             0x0002  /* Route entry for default gateway */
494 #define IRE_LOCAL               0x0004  /* Route entry for local address */
495 #define IRE_LOOPBACK            0x0008  /* Route entry for loopback address */
496 #define IRE_PREFIX              0x0010  /* Route entry for prefix routes */
497 #ifndef _KERNEL
498 /* Keep so user-level still compiles */
499 #define IRE_CACHE               0x0020  /* Cached Route entry */
500 #endif
501 #define IRE_IF_NORESOLVER       0x0040  /* Route entry for local interface */
502                                         /* net without any address mapping. */
503 #define IRE_IF_RESOLVER         0x0080  /* Route entry for local interface */
504                                         /* net with resolver. */
505 #define IRE_HOST                0x0100  /* Host route entry */
506 /* Keep so user-level still compiles */
507 #define IRE_HOST_REDIRECT       0x0200  /* only used for T_SVR4_OPTMGMT_REQ */
508 #define IRE_IF_CLONE            0x0400  /* Per host clone of IRE_IF */
509 #define IRE_MULTICAST           0x0800  /* Special - not in table */
510 #define IRE_NOROUTE             0x1000  /* Special - not in table */

512 #define IRE_INTERFACE           (IRE_IF_NORESOLVER | IRE_IF_RESOLVER)

514 #define IRE_IF_ALL              (IRE_IF_NORESOLVER | IRE_IF_RESOLVER | \
515                                 IRE_IF_CLONE)
516 #define IRE_OFFSUBNET           (IRE_DEFAULT | IRE_PREFIX | IRE_HOST)
517 #define IRE_OFFLINK             IRE_OFFSUBNET
518 /*
519  * Note that we view IRE_NOROUTE as ONLINK since we can "send" to them without
520  * going through a router; the result of sending will be an error/icmp error.
521  */
522 #define IRE_ONLINK              (IRE_IF_ALL|IRE_LOCAL|IRE_LOOPBACK| \
523                                 IRE_BROADCAST|IRE_MULTICAST|IRE_NOROUTE)
```

```
525 /* Arguments to ire_flush_cache() */
526 #define IRE_FLUSH_DELETE        0
527 #define IRE_FLUSH_ADD           1
528 #define IRE_FLUSH_GWCHANGE      2

530 /*
531  * Flags to ire_route_recursive
532  */
533 #define IRR_NONE                0
534 #define IRR_ALLOCATE            1       /* OK to allocate IRE_IF_CLONE */
535 #define IRR_INCOMPLETE          2       /* OK to return incomplete chain */

537 /*
538  * Open/close synchronization flags.
539  * These are kept in a separate field in the conn and the synchronization
540  * depends on the atomic 32 bit access to that field.
541  */
542 #define CONN_CLOSING            0x01    /* ip_close waiting for ip_wsrv */
543 #define CONN_CONDEMNED          0x02    /* conn is closing, no more refs */
544 #define CONN_INCIPIENT          0x04    /* conn not yet visible, no refs */
545 #define CONN_QUIESCED           0x08    /* conn is now quiescent */
546 #define CONN_UPDATE_ILL         0x10    /* conn_update_ill in progress */

548 /*
549  * Flags for dce_flags field. Specifies which information has been set.
550  * dce_ident is always present, but the other ones are identified by the flags.
551  */
552 #define DCEF_DEFAULT            0x0001  /* Default DCE - no pmtu or uinfo */
553 #define DCEF_PMTU               0x0002  /* Different than interface MTU */
554 #define DCEF_UINFO              0x0004  /* dce_uinfo set */
555 #define DCEF_TOO_SMALL_PMTU     0x0008  /* Smaller than IPv4/IPv6 MIN */

557 #ifdef _KERNEL
558 /*
559  * Extra structures need for per-src-addr filtering (IGMPv3/MLDv2)
560  */
561 #define MAX_FILTER_SIZE 64

563 typedef struct slist_s {
564         int             sl_numsrc;
565         in6_addr_t      sl_addr[MAX_FILTER_SIZE];
566 } slist_t;

568 /*
569  * Following struct is used to maintain retransmission state for
570  * a multicast group.  One rtx_state_t struct is an in-line field
571  * of the ilm_t struct; the slist_ts in the rtx_state_t struct are
572  * alloc'd as needed.
573  */
574 typedef struct rtx_state_s {
575         uint_t          rtx_timer;      /* retrans timer */
576         int             rtx_cnt;        /* retrans count */
577         int             rtx_fmode_cnt;  /* retrans count for fmode change */
578         slist_t         *rtx_allow;
579         slist_t         *rtx_block;
580 } rtx_state_t;

582 /*
583  * Used to construct list of multicast address records that will be
584  * sent in a single listener report.
585  */
586 typedef struct mrec_s {
587         struct mrec_s   *mrec_next;
588         uint8_t         mrec_type;
589         uint8_t         mrec_auxlen;    /* currently unused */
```

```
590         in6_addr_t      mrec_group;
591         slist_t         mrec_srcs;
592 } mrec_t;

594 /* Group membership list per upper conn */

596 /*
597  * We record the multicast information from the socket option in
598  * ilg_ifaddr/ilg_ifindex. This allows rejoining the group in the case when
599  * the ifaddr (or ifindex) disappears and later reappears, potentially on
600  * a different ill. The IPv6 multicast socket options and ioctls all specify
601  * the interface using an ifindex. For IPv4 some socket options/ioctls use
602  * the interface address and others use the index. We record here the method
603  * that was actually used (and leave the other of ilg_ifaddr or ilg_ifindex)
604  * at zero so that we can rejoin the way the application intended.
605  *
606  * We track the ill on which we will or already have joined an ilm using
607  * ilg_ill. When we have succeeded joining the ilm and have a refhold on it
608  * then we set ilg_ilm. Thus intentionally there is a window where ilg_ill is
609  * set and ilg_ilm is not set. This allows clearing ilg_ill as a signal that
610  * the ill is being unplumbed and the ilm should be discarded.
611  *
612  * ilg records the state of multicast memberships of a socket end point.
613  * ilm records the state of multicast memberships with the driver and is
614  * maintained per interface.
615  *
616  * The ilg state is protected by conn_ilg_lock.
617  * The ilg will not be freed until ilg_refcnt drops to zero.
618  */
619 typedef struct ilg_s {
620         struct ilg_s    *ilg_next;
621         struct ilg_s    **ilg_ptpn;
622         struct conn_s   *ilg_connp;     /* Back pointer to get lock */
623         in6_addr_t      ilg_v6group;
624         ipaddr_t        ilg_ifaddr;     /* For some IPv4 cases */
625         uint_t          ilg_ifindex;    /* IPv6 and some other IPv4 cases */
626         struct ill_s    *ilg_ill;       /* Where ilm is joined. No refhold */
627         struct ilm_s    *ilg_ilm;       /* With ilm_refhold */
628         uint_t          ilg_refcnt;
629         mcast_record_t  ilg_fmode;      /* MODE_IS_INCLUDE/MODE_IS_EXCLUDE */
630         slist_t         *ilg_filter;
631         boolean_t       ilg_condemned;  /* Conceptually deleted */
632 } ilg_t;

634 /*
635  * Multicast address list entry for ill.
636  * ilm_ill is used by IPv4 and IPv6
637  *
638  * The ilm state (and other multicast state on the ill) is protected by
639  * ill_mcast_lock. Operations that change state on both an ilg and ilm
640  * in addition use ill_mcast_serializer to ensure that we can't have
641  * interleaving between e.g., add and delete operations for the same conn_t,
642  * group, and ill. The ill_mcast_serializer is also used to ensure that
643  * multicast group joins do not occur on an interface that is in the process
644  * of joining an IPMP group.
645  *
646  * The comment below (and for other netstack_t references) refers
647  * to the fact that we only do netstack_hold in particular cases,
648  * such as the references from open endpoints (ill_t and conn_t's
649  * pointers). Internally within IP we rely on IP's ability to cleanup e.g.
650  * ire_t's when an ill goes away.
651  */
652 typedef struct ilm_s {
653         in6_addr_t      ilm_v6addr;
654         int             ilm_refcnt;
655         uint_t          ilm_timer;      /* IGMP/MLD query resp timer, in msec */
```

```
656         struct ilm_s    *ilm_next;      /* Linked list for each ill */
657         uint_t          ilm_state;      /* state of the membership */
658         struct ill_s    *ilm_ill;       /* Back pointer to ill - ill_ilm_cnt */
659         zoneid_t        ilm_zoneid;
660         int             ilm_no_ilg_cnt; /* number of joins w/ no ilg */
661         mcast_record_t  ilm_fmode;      /* MODE_IS_INCLUDE/MODE_IS_EXCLUDE */
662         slist_t         *ilm_filter;    /* source filter list */
663         slist_t         *ilm_pendsrcs;  /* relevant src addrs for pending req */
664         rtx_state_t     ilm_rtx;        /* SCR retransmission state */
665         ipaddr_t        ilm_ifaddr;     /* For IPv4 netstat */
666         ip_stack_t      *ilm_ipst;      /* Does not have a netstack_hold */
667 } ilm_t;

669 #define ilm_addr        V4_PART_OF_V6(ilm_v6addr)


671 /*
672  * Soft reference to an IPsec SA.
673  *
674  * On relative terms, conn's can be persistent (living as long as the
675  * processes which create them), while SA's are ephemeral (dying when
676  * they hit their time-based or byte-based lifetimes).
677  *
678  * We could hold a hard reference to an SA from an ipsec_latch_t,
679  * but this would cause expired SA's to linger for a potentially
680  * unbounded time.
681  *
682  * Instead, we remember the hash bucket number and bucket generation
683  * in addition to the pointer.  The bucket generation is incremented on
684  * each deletion.
685  */
686 typedef struct ipsa_ref_s
687 {
688         struct ipsa_s   *ipsr_sa;
689         struct isaf_s   *ipsr_bucket;
690         uint64_t        ipsr_gen;
691 } ipsa_ref_t;

693 /*
694  * IPsec "latching" state.
695  *
696  * In the presence of IPsec policy, fully-bound conn's bind a connection
697  * to more than just the 5-tuple, but also a specific IPsec action and
698  * identity-pair.
699  * The identity pair is accessed from both the receive and transmit side
700  * hence it is maintained in the ipsec_latch_t structure. conn_latch and
701  * ixa_ipsec_latch points to it.
702  * The policy and actions are stored in conn_latch_in_policy and
703  * conn_latch_in_action for the inbound side, and in ixa_ipsec_policy and
704  * ixa_ipsec_action for the transmit side.
705  *
706  * As an optimization, we also cache soft references to IPsec SA's in
707  * ip_xmit_attr_t so that we can fast-path around most of the work needed for
708  * outbound IPsec SA selection.
709  */
710 typedef struct ipsec_latch_s
711 {
712         kmutex_t        ipl_lock;
713         uint32_t        ipl_refcnt;

715         struct ipsid_s *ipl_local_cid;
716         struct ipsid_s *ipl_remote_cid;
717         unsigned int
718                         ipl_ids_latched : 1,

720                         ipl_pad_to_bit_31 : 31;
721 } ipsec_latch_t;
```

```
723 #define IPLATCH_REFHOLD(ipl) { \
724         atomic_add_32(&(ipl)->ipl_refcnt, 1);           \
725         ASSERT((ipl)->ipl_refcnt != 0);                 \
726 }

728 #define IPLATCH_REFRELE(ipl) {                          \
729         ASSERT((ipl)->ipl_refcnt != 0);                         \
730         membar_exit();                                          \
731         if (atomic_add_32_nv(&(ipl)->ipl_refcnt, -1) == 0)      \
732                 iplatch_free(ipl);                              \
733 }

735 /*
736  * peer identity structure.
737  */
738 typedef struct conn_s conn_t;

740 /*
741  * This is used to match an inbound/outbound datagram with policy.
742  */
743 typedef struct ipsec_selector {
744         in6_addr_t      ips_local_addr_v6;
745         in6_addr_t      ips_remote_addr_v6;
746         uint16_t        ips_local_port;
747         uint16_t        ips_remote_port;
748         uint8_t         ips_icmp_type;
749         uint8_t         ips_icmp_code;
750         uint8_t         ips_protocol;
751         uint8_t         ips_isv4 : 1,
752                         ips_is_icmp_inv_acq: 1;
753 } ipsec_selector_t;

755 /*
756  * Note that we put v4 addresses in the *first* 32-bit word of the
757  * selector rather than the last to simplify the prefix match/mask code
758  * in spd.c
759  */
760 #define ips_local_addr_v4 ips_local_addr_v6.s6_addr32[0]
761 #define ips_remote_addr_v4 ips_remote_addr_v6.s6_addr32[0]

763 /* Values used in IP by IPSEC Code */
764 #define         IPSEC_OUTBOUND          B_TRUE
765 #define         IPSEC_INBOUND           B_FALSE

767 /*
768  * There are two variants in policy failures. The packet may come in
769  * secure when not needed (IPSEC_POLICY_???_NOT_NEEDED) or it may not
770  * have the desired level of protection (IPSEC_POLICY_MISMATCH).
771  */
772 #define IPSEC_POLICY_NOT_NEEDED         0
773 #define IPSEC_POLICY_MISMATCH           1
774 #define IPSEC_POLICY_AUTH_NOT_NEEDED    2
775 #define IPSEC_POLICY_ENCR_NOT_NEEDED    3
776 #define IPSEC_POLICY_SE_NOT_NEEDED      4
777 #define IPSEC_POLICY_MAX                5       /* Always max + 1. */

779 /*
780  * Check with IPSEC inbound policy if
781  *
782  * 1) per-socket policy is present - indicated by conn_in_enforce_policy.
783  * 2) Or if we have not cached policy on the conn and the global policy is
784  *    non-empty.
785  */
786 #define CONN_INBOUND_POLICY_PRESENT(connp, ipss)        \
787         ((connp)->conn_in_enforce_policy ||             \
```

```
788             (!((connp)->conn_policy_cached) &&                    \
789             (ipss)->ipsec_inbound_v4_policy_present))

791 #define CONN_INBOUND_POLICY_PRESENT_V6(connp, ipss)     \
792             ((connp)->conn_in_enforce_policy ||          \
793             (!(connp)->conn_policy_cached &&             \
794             (ipss)->ipsec_inbound_v6_policy_present))

796 #define CONN_OUTBOUND_POLICY_PRESENT(connp, ipss)       \
797             ((connp)->conn_out_enforce_policy ||         \
798             (!((connp)->conn_policy_cached) &&           \
799             (ipss)->ipsec_outbound_v4_policy_present))

801 #define CONN_OUTBOUND_POLICY_PRESENT_V6(connp, ipss)    \
802             ((connp)->conn_out_enforce_policy ||         \
803             (!(connp)->conn_policy_cached &&            \
804             (ipss)->ipsec_outbound_v6_policy_present))
806 /*
807  * Information cached in IRE for upper layer protocol (ULP).
808  */
809 typedef struct iulp_s {
810         boolean_t       iulp_set;       /* Is any metric set? */
811         uint32_t        iulp_ssthresh;  /* Slow start threshold (TCP). */
812         clock_t         iulp_rtt;       /* Guestimate in millisecs. */
813         clock_t         iulp_rtt_sd;    /* Cached value of RTT variance. */
814         uint32_t        iulp_spipe;     /* Send pipe size. */
815         uint32_t        iulp_rpipe;     /* Receive pipe size. */
816         uint32_t        iulp_rtomax;    /* Max round trip timeout. */
817         uint32_t        iulp_sack;      /* Use SACK option (TCP)? */
818         uint32_t        iulp_mtu;       /* Setable with routing sockets */

820         uint32_t
821                 iulp_tstamp_ok : 1,     /* Use timestamp option (TCP)? */
822                 iulp_wscale_ok : 1,     /* Use window scale option (TCP)? */
823                 iulp_ecn_ok : 1,        /* Enable ECN (for TCP)? */
824                 iulp_pmtud_ok : 1,      /* Enable PMTUd? */

826                 /* These three are passed out by ip_set_destination */
827                 iulp_localnet: 1,       /* IRE_ONLINK */
828                 iulp_loopback: 1,       /* IRE_LOOPBACK */
829                 iulp_local: 1,          /* IRE_LOCAL */

831                 iulp_not_used : 25;
832 } iulp_t;

834 /*
835  * The conn drain list structure (idl_t), protected by idl_lock.  Each conn_t
836  * inserted in the list points back at this idl_t using conn_idl, and is
837  * chained by conn_drain_next and conn_drain_prev, which are also protected by
838  * idl_lock.  When flow control is relieved, either ip_wsrv() (STREAMS) or
839  * ill_flow_enable() (non-STREAMS) will call conn_drain().
840  *
841  * The conn drain list, idl_t, itself is part of tx cookie list structure.
842  * A tx cookie list points to a blocked Tx ring and contains the list of
843  * all conn's that are blocked due to the flow-controlled Tx ring (via
844  * the idl drain list). Note that a link can have multiple Tx rings. The
845  * drain list will store the conn's blocked due to Tx ring being flow
846  * controlled.
847  */

849 typedef uintptr_t ip_mac_tx_cookie_t;
850 typedef struct idl_s idl_t;
851 typedef struct idl_tx_list_s idl_tx_list_t;

853 struct idl_tx_list_s {
```

```
854         ip_mac_tx_cookie_t      txl_cookie;
855         kmutex_t                txl_lock;       /* Lock for this list */
856         idl_t                   *txl_drain_list;
857         int                     txl_drain_index;
858 };

860 struct idl_s {
861         conn_t          *idl_conn;              /* Head of drain list */
862         kmutex_t        idl_lock;               /* Lock for this list */
863         idl_tx_list_t   *idl_itl;
864 };

866 /*
867  * Interface route structure which holds the necessary information to recreate
868  * routes that are tied to an interface i.e. have ire_ill set.
869  *
870  * These routes which were initially created via a routing socket or via the
871  * SIOCADDRT ioctl may be gateway routes (RTF_GATEWAY being set) or may be
872  * traditional interface routes.  When an ill comes back up after being
873  * down, this information will be used to recreate the routes.  These
874  * are part of an mblk_t chain that hangs off of the ILL (ill_saved_ire_mp).
875  */
876 typedef struct ifrt_s {
877         ushort_t        ifrt_type;              /* Type of IRE */
878         in6_addr_t      ifrt_v6addr;            /* Address IRE represents. */
879         in6_addr_t      ifrt_v6gateway_addr;    /* Gateway if IRE_OFFLINK */
880         in6_addr_t      ifrt_v6setsrc_addr;     /* Src addr if RTF_SETSRC */
881         in6_addr_t      ifrt_v6mask;            /* Mask for matching IRE. */
882         uint32_t        ifrt_flags;             /* flags related to route */
883         iulp_t          ifrt_metrics;           /* Routing socket metrics */
884         zoneid_t        ifrt_zoneid;            /* zoneid for route */
885 } ifrt_t;

887 #define ifrt_addr               V4_PART_OF_V6(ifrt_v6addr)
888 #define ifrt_gateway_addr       V4_PART_OF_V6(ifrt_v6gateway_addr)
889 #define ifrt_mask               V4_PART_OF_V6(ifrt_v6mask)
890 #define ifrt_setsrc_addr        V4_PART_OF_V6(ifrt_v6setsrc_addr)

892 /* Number of IP addresses that can be hosted on a physical interface */
893 #define MAX_ADDRS_PER_IF        8192
894 /*
895  * Number of Source addresses to be considered for source address
896  * selection. Used by ipif_select_source_v4/v6.
897  */
898 #define MAX_IPIF_SELECT_SOURCE  50

900 #ifdef IP_DEBUG
901 /*
902  * Trace refholds and refreles for debugging.
903  */
904 #define TR_STACK_DEPTH 14
905 typedef struct tr_buf_s {
906         int     tr_depth;
907         clock_t tr_time;
908         pc_t    tr_stack[TR_STACK_DEPTH];
909 } tr_buf_t;

911 typedef struct th_trace_s {
912         int             th_refcnt;
913         uint_t          th_trace_lastref;
914         kthread_t       *th_id;
915 #define TR_BUF_MAX      38
916         tr_buf_t        th_trbuf[TR_BUF_MAX];
917 } th_trace_t;

919 typedef struct th_hash_s {
```

```
 920          list_node_t     thh_link;
 921          mod_hash_t      *thh_hash;
 922          ip_stack_t      *thh_ipst;
 923 } th_hash_t;
 924 #endif

 926 /* The following are ipif_state_flags */
 927 #define IPIF_CONDEMNED          0x1     /* The ipif is being removed */
 928 #define IPIF_CHANGING           0x2     /* A critcal ipif field is changing */
 929 #define IPIF_SET_LINKLOCAL      0x10    /* transient flag during bringup */

 931 /* IP interface structure, one per local address */
 932 typedef struct ipif_s {
 933          struct  ipif_s *ipif_next;
 934          struct  ill_s  *ipif_ill;     /* Back pointer to our ill */
 935          int     ipif_id;              /* Logical unit number */
 936          in6_addr_t ipif_v6lcl_addr;   /* Local IP address for this if. */
 937          in6_addr_t ipif_v6subnet;     /* Subnet prefix for this if. */
 938          in6_addr_t ipif_v6net_mask;   /* Net mask for this interface. */
 939          in6_addr_t ipif_v6brd_addr;   /* Broadcast addr for this interface. */
 940          in6_addr_t ipif_v6pp_dst_addr; /* Point-to-point dest address. */
 941          uint64_t ipif_flags;          /* Interface flags. */
 942          uint_t  ipif_ire_type;        /* IRE_LOCAL or IRE_LOOPBACK */

 944          /*
 945           * The packet count in the ipif contain the sum of the
 946           * packet counts in dead IRE_LOCAL/LOOPBACK for this ipif.
 947           */
 948          uint_t  ipif_ib_pkt_count;    /* Inbound packets for our dead IREs */

 950          /* Exclusive bit fields, protected by ipsq_t */
 951          unsigned int
 952                  ipif_was_up : 1,      /* ipif was up before */
 953                  ipif_addr_ready : 1,  /* DAD is done */
 954                  ipif_was_dup : 1,     /* DAD had failed */
 955                  ipif_added_nce : 1,   /* nce added for local address */

 957                  ipif_pad_to_31 : 28;

 959          ilm_t  *ipif_allhosts_ilm;    /* For all-nodes join */
 960          ilm_t  *ipif_solmulti_ilm;    /* For IPv6 solicited multicast join */

 962          uint_t  ipif_seqid;           /* unique index across all ills */
 963          uint_t  ipif_state_flags;     /* See IPIF_* flag defs above */
 964          uint_t  ipif_refcnt;          /* active consistent reader cnt */

 966          zoneid_t ipif_zoneid;         /* zone ID number */
 967          timeout_id_t ipif_recovery_id; /* Timer for DAD recovery */
 968          boolean_t ipif_trace_disable;  /* True when alloc fails */
 969          /*
 970           * For an IPMP interface, ipif_bound_ill tracks the ill whose hardware
 971           * information this ipif is associated with via ARP/NDP.  We can use
 972           * an ill pointer (rather than an index) because only ills that are
 973           * part of a group will be pointed to, and an ill cannot disappear
 974           * while it's in a group.
 975           */
 976          struct ill_s   *ipif_bound_ill;
 977          struct ipif_s  *ipif_bound_next; /* bound ipif chain */
 978          boolean_t       ipif_bound;    /* B_TRUE if we successfully bound */

 980          struct ire_s   *ipif_ire_local; /* Our IRE_LOCAL or LOOPBACK */
 981          struct ire_s   *ipif_ire_if;    /* Our IRE_INTERFACE */
 982 } ipif_t;

 984 /*
 985  * The following table lists the protection levels of the various members
```

```
 986  * of the ipif_t. The following notation is used.
 987  *
 988  * Write once - Written to only once at the time of bringing up
 989  * the interface and can be safely read after the bringup without any lock.
 990  *
 991  * ipsq - Need to execute in the ipsq to perform the indicated access.
 992  *
 993  * ill_lock - Need to hold this mutex to perform the indicated access.
 994  *
 995  * ill_g_lock - Need to hold this rw lock as reader/writer for read access or
 996  * write access respectively.
 997  *
 998  * down ill - Written to only when the ill is down (i.e all ipifs are down)
 999  * up ill - Read only when the ill is up (i.e. at least 1 ipif is up)
1000  *
1001  *              Table of ipif_t members and their protection
1002  *
1003  * ipif_next            ipsq + ill_lock +       ipsq OR ill_lock OR
1004  *                      ill_g_lock              ill_g_lock
1005  * ipif_ill             ipsq + down ipif        write once
1006  * ipif_id              ipsq + down ipif        write once
1007  * ipif_v6lcl_addr      ipsq + down ipif        up ipif
1008  * ipif_v6subnet        ipsq + down ipif        up ipif
1009  * ipif_v6net_mask      ipsq + down ipif        up ipif
1010  *
1011  * ipif_v6brd_addr
1012  * ipif_v6pp_dst_addr
1013  * ipif_flags           ill_lock                ill_lock
1014  * ipif_ire_type        ipsq + down ill         up ill
1015  *
1016  * ipif_ib_pkt_count    Approx
1017  *
1018  * bit fields           ill_lock                ill_lock
1019  *
1020  * ipif_allhosts_ilm    ipsq                    ipsq
1021  * ipif_solmulti_ilm    ipsq                    ipsq
1022  *
1023  * ipif_seqid           ipsq                    Write once
1024  *
1025  * ipif_state_flags     ill_lock                ill_lock
1026  * ipif_refcnt          ill_lock                ill_lock
1027  * ipif_bound_ill       ipsq + ipmp_lock        ipsq OR ipmp_lock
1028  * ipif_bound_next      ipsq                    ipsq
1029  * ipif_bound           ipsq                    ipsq
1030  *
1031  * ipif_ire_local       ipsq + ips_ill_g_lock   ipsq OR ips_ill_g_lock
1032  * ipif_ire_if          ipsq + ips_ill_g_lock   ipsq OR ips_ill_g_lock
1033  */

1035 /*
1036  * Return values from ip_laddr_verify_{v4,v6}
1037  */
1038 typedef enum { IPVL_UNICAST_UP, IPVL_UNICAST_DOWN, IPVL_MCAST, IPVL_BCAST,
1039              IPVL_BAD} ip_laddr_t;


1042 #define IP_TR_HASH(tid) ((((uintptr_t)tid) >> 6) & (IP_TR_HASH_MAX - 1))

1044 #ifdef DEBUG
1045 #define IPIF_TRACE_REF(ipif)    ipif_trace_ref(ipif)
1046 #define ILL_TRACE_REF(ill)      ill_trace_ref(ill)
1047 #define IPIF_UNTRACE_REF(ipif)  ipif_untrace_ref(ipif)
1048 #define ILL_UNTRACE_REF(ill)    ill_untrace_ref(ill)
1049 #else
1050 #define IPIF_TRACE_REF(ipif)
1051 #define ILL_TRACE_REF(ill)
```

```
1052 #define IPIF_UNTRACE_REF(ipif)
1053 #define ILL_UNTRACE_REF(ill)
1054 #endif

1056 /* IPv4 compatibility macros */
1057 #define ipif_lcl_addr           V4_PART_OF_V6(ipif_v6lcl_addr)
1058 #define ipif_subnet             V4_PART_OF_V6(ipif_v6subnet)
1059 #define ipif_net_mask           V4_PART_OF_V6(ipif_v6net_mask)
1060 #define ipif_brd_addr           V4_PART_OF_V6(ipif_v6brd_addr)
1061 #define ipif_pp_dst_addr        V4_PART_OF_V6(ipif_v6pp_dst_addr)

1063 /* Macros for easy backreferences to the ill. */
1064 #define ipif_isv6               ipif_ill->ill_isv6

1066 #define SIOCLIFADDR_NDX 112     /* ndx of SIOCLIFADDR in the ndx ioctl table */

1068 /*
1069  * mode value for ip_ioctl_finish for finishing an ioctl
1070  */
1071 #define CONN_CLOSE      1               /* No mi_copy */
1072 #define COPYOUT         2               /* do an mi_copyout if needed */
1073 #define NO_COPYOUT      3               /* do an mi_copy_done */
1074 #define IPI2MODE(ipi)   ((ipi)->ipi_flags & IPI_GET_CMD ? COPYOUT : NO_COPYOUT)

1076 /*
1077  * The IP-MT design revolves around the serialization objects ipsq_t (IPSQ)
1078  * and ipxop_t (exclusive operation or "xop").  Becoming "writer" on an IPSQ
1079  * ensures that no other threads can become "writer" on any IPSQs sharing that
1080  * IPSQ's xop until the writer thread is done.
1081  *
1082  * Each phyint points to one IPSQ that remains fixed over the phyint's life.
1083  * Each IPSQ points to one xop that can change over the IPSQ's life.  If a
1084  * phyint is *not* in an IPMP group, then its IPSQ will refer to the IPSQ's
1085  * "own" xop (ipsq_ownxop).  If a phyint *is* part of an IPMP group, then its
1086  * IPSQ will refer to the "group" xop, which is shorthand for the xop of the
1087  * IPSQ of the IPMP meta-interface's phyint.  Thus, all phyints that are part
1088  * of the same IPMP group will have their IPSQ's point to the group xop, and
1089  * thus becoming "writer" on any phyint in the group will prevent any other
1090  * writer on any other phyint in the group.  All IPSQs sharing the same xop
1091  * are chained together through ipsq_next (in the degenerate common case,
1092  * ipsq_next simply refers to itself).  Note that the group xop is guaranteed
1093  * to exist at least as long as there are members in the group, since the IPMP
1094  * meta-interface can only be destroyed if the group is empty.
1095  *
1096  * Incoming exclusive operation requests are enqueued on the IPSQ they arrived
1097  * on rather than the xop.  This makes switching xop's (as would happen when a
1098  * phyint leaves an IPMP group) simple, because after the phyint leaves the
1099  * group, any operations enqueued on its IPSQ can be safely processed with
1100  * respect to its new xop, and any operations enqueued on the IPSQs of its
1101  * former group can be processed with respect to their existing group xop.
1102  * Even so, switching xops is a subtle dance; see ipsq_dq() for details.
1103  *
1104  * An IPSQ's "own" xop is embedded within the IPSQ itself since they have have
1105  * identical lifetimes, and because doing so simplifies pointer management.
1106  * While each phyint and IPSQ point to each other, it is not possible to free
1107  * the IPSQ when the phyint is freed, since we may still *inside* the IPSQ
1108  * when the phyint is being freed.  Thus, ipsq_phyint is set to NULL when the
1109  * phyint is freed, and the IPSQ free is later done in ipsq_exit().
1110  *
1111  * ipsq_t synchronization:        read                    write
1112  *
1113  *      ipsq_xopq_mphead        ipx_lock                ipx_lock
1114  *      ipsq_xopq_mptail        ipx_lock                ipx_lock
1115  *      ipsq_xop_switch_mp      ipsq_lock               ipsq_lock
1116  *      ipsq_phyint             write once              write once
1117  *      ipsq_next               RW_READER ill_g_lock    RW_WRITER ill_g_lock
```

```
1118  *      ipsq_xop                ipsq_lock or ipsq       ipsq_lock + ipsq
1119  *      ipsq_swxop              ipsq                    ipsq
1120  *      ipsq_ownxop             see ipxop_t             see ipxop_t
1121  *      ipsq_ipst               write once              write once
1122  *
1123  * ipxop_t synchronization:       read                    write
1124  *
1125  *      ipx_writer              ipx_lock                ipx_lock
1126  *      ipx_xop_queued          ipx_lock                ipx_lock
1127  *      ipx_mphead              ipx_lock                ipx_lock
1128  *      ipx_mptail              ipx_lock                ipx_lock
1129  *      ipx_ipsq                write once              write once
1130  *      ips_ipsq_queued         ipx_lock                ipx_lock
1131  *      ipx_waitfor             ipsq or ipx_lock        ipsq + ipx_lock
1132  *      ipx_reentry_cnt         ipsq or ipx_lock        ipsq + ipx_lock
1133  *      ipx_current_done        ipsq                    ipsq
1134  *      ipx_current_ioctl       ipsq                    ipsq
1135  *      ipx_current_ipif        ipsq or ipx_lock        ipsq + ipx_lock
1136  *      ipx_pending_ipif        ipsq or ipx_lock        ipsq + ipx_lock
1137  *      ipx_pending_mp          ipsq or ipx_lock        ipsq + ipx_lock
1138  *      ipx_forced              ipsq                    ipsq
1139  *      ipx_depth               ipsq                    ipsq
1140  *      ipx_stack               ipsq                    ipsq
1141  */
1142 typedef struct ipxop_s {
1143         kmutex_t        ipx_lock;       /* see above */
1144         kthread_t       *ipx_writer;    /* current owner */
1145         mblk_t          *ipx_mphead;    /* messages tied to this op */
1146         mblk_t          *ipx_mptail;
1147         struct ipsq_s   *ipx_ipsq;      /* associated ipsq */
1148         boolean_t       ipx_ipsq_queued; /* ipsq using xop has queued op */
1149         int             ipx_waitfor;    /* waiting; values encoded below */
1150         int             ipx_reentry_cnt;
1151         boolean_t       ipx_current_done;  /* is the current operation done? */
1152         int             ipx_current_ioctl; /* current ioctl, or 0 if no ioctl */
1153         ipif_t          *ipx_current_ipif; /* ipif for current op */
1154         ipif_t          *ipx_pending_ipif; /* ipif for ipx_pending_mp */
1155         mblk_t          *ipx_pending_mp;   /* current ioctl mp while waiting */
1156         boolean_t       ipx_forced;                     /* debugging aid */
1157 #ifdef DEBUG
1158         int             ipx_depth;                      /* debugging aid */
1159 #define IPX_STACK_DEPTH 15
1160         pc_t            ipx_stack[IPX_STACK_DEPTH];      /* debugging aid */
1161 #endif
1162 } ipxop_t;

1164 typedef struct ipsq_s {
1165         kmutex_t ipsq_lock;             /* see above */
1166         mblk_t  *ipsq_switch_mp;        /* op to handle right after switch */
1167         mblk_t  *ipsq_xopq_mphead;      /* list of excl ops (mostly ioctls) */
1168         mblk_t  *ipsq_xopq_mptail;
1169         struct phyint  *ipsq_phyint;    /* associated phyint */
1170         struct ipsq_s  *ipsq_next;      /* next ipsq sharing ipsq_xop */
1171         struct ipxop_s *ipsq_xop;       /* current xop synchronization info */
1172         struct ipxop_s *ipsq_swxop;     /* switch xop to on ipsq_exit() */
1173         struct ipxop_s ipsq_ownxop;     /* our own xop (may not be in-use) */
1174         ip_stack_t     *ipsq_ipst;      /* does not have a netstack_hold */
1175 } ipsq_t;

1177 /*
1178  * ipx_waitfor values:
1179  */
1180 enum {
1181         IPIF_DOWN = 1,  /* ipif_down() waiting for refcnts to drop */
1182         ILL_DOWN,       /* ill_down() waiting for refcnts to drop */
1183         IPIF_FREE,      /* ipif_free() waiting for refcnts to drop */
```

```
1184         ILL_FREE          /* ill unplumb waiting for refcnts to drop */
1185 };

1187 /* Operation types for ipsq_try_enter() */
1188 #define CUR_OP 0          /* request writer within current operation */
1189 #define NEW_OP 1          /* request writer for a new operation */
1190 #define SWITCH_OP 2       /* request writer once IPSQ XOP switches */

1192 /*
1193  * Kstats tracked on each IPMP meta-interface.  Order here must match
1194  * ipmp_kstats[] in ip/ipmp.c.
1195  */
1196 enum {
1197         IPMP_KSTAT_OBYTES,      IPMP_KSTAT_OBYTES64,    IPMP_KSTAT_RBYTES,
1198         IPMP_KSTAT_RBYTES64,    IPMP_KSTAT_OPACKETS,    IPMP_KSTAT_OPACKETS64,
1199         IPMP_KSTAT_OERRORS,     IPMP_KSTAT_IPACKETS,    IPMP_KSTAT_IPACKETS64,
1200         IPMP_KSTAT_IERRORS,     IPMP_KSTAT_MULTIRCV,    IPMP_KSTAT_MULTIXMT,
1201         IPMP_KSTAT_BRDCSTRCV,   IPMP_KSTAT_BRDCSTXMT,   IPMP_KSTAT_LINK_UP,
1202         IPMP_KSTAT_MAX          /* keep last */
1203 };

1205 /*
1206  * phyint represents state that is common to both IPv4 and IPv6 interfaces.
1207  * There is a separate ill_t representing IPv4 and IPv6 which has a
1208  * backpointer to the phyint structure for accessing common state.
1209  */
1210 typedef struct phyint {
1211         struct ill_s    *phyint_illv4;
1212         struct ill_s    *phyint_illv6;
1213         uint_t          phyint_ifindex;         /* SIOCSLIFINDEX */
1214         uint64_t        phyint_flags;
1215         avl_node_t      phyint_avl_by_index;    /* avl tree by index */
1216         avl_node_t      phyint_avl_by_name;     /* avl tree by name */
1217         kmutex_t        phyint_lock;
1218         struct ipsq_s   *phyint_ipsq;           /* back pointer to ipsq */
1219         struct ipmp_grp_s *phyint_grp;          /* associated IPMP group */
1220         char            phyint_name[LIFNAMSIZ]; /* physical interface name */
1221         uint64_t        phyint_kstats0[IPMP_KSTAT_MAX]; /* baseline kstats */
1222 } phyint_t;

1224 #define CACHE_ALIGN_SIZE 64
1225 #define CACHE_ALIGN(align_struct)       P2ROUNDUP(sizeof (struct align_struct),\
1226                                                 CACHE_ALIGN_SIZE)
1227 struct _phyint_list_s_ {
1228         avl_tree_t      phyint_list_avl_by_index;       /* avl tree by index */
1229         avl_tree_t      phyint_list_avl_by_name;        /* avl tree by name */
1230 };

1232 typedef union phyint_list_u {
1233         struct  _phyint_list_s_ phyint_list_s;
1234         char    phyint_list_filler[CACHE_ALIGN(_phyint_list_s_)];
1235 } phyint_list_t;

1237 #define phyint_list_avl_by_index        phyint_list_s.phyint_list_avl_by_index
1238 #define phyint_list_avl_by_name         phyint_list_s.phyint_list_avl_by_name

1240 /*
1241  * Fragmentation hash bucket
1242  */
1243 typedef struct ipfb_s {
1244         struct ipf_s    *ipfb_ipf;      /* List of ... */
1245         size_t          ipfb_count;     /* Count of bytes used by frag(s) */
1246         kmutex_t        ipfb_lock;      /* Protect all ipf in list */
1247         uint_t          ipfb_frag_pkts; /* num of distinct fragmented pkts */
1248 } ipfb_t;
```

```
1250 /*
1251  * IRE bucket structure. Usually there is an array of such structures,
1252  * each pointing to a linked list of ires. irb_refcnt counts the number
1253  * of walkers of a given hash bucket. Usually the reference count is
1254  * bumped up if the walker wants no IRES to be DELETED while walking the
1255  * list. Bumping up does not PREVENT ADDITION. This allows walking a given
1256  * hash bucket without stumbling up on a free pointer.
1257  *
1258  * irb_t structures in ip_ftable are dynamically allocated and freed.
1259  * In order to identify the irb_t structures that can be safely kmem_free'd
1260  * we need to ensure that
1261  *  - the irb_refcnt is quiescent, indicating no other walkers,
1262  *  - no other threads or ire's are holding references to the irb,
1263  *      i.e., irb_nire == 0,
1264  *  - there are no active ire's in the bucket, i.e., irb_ire_cnt == 0
1265  */
1266 typedef struct irb {
1267         struct ire_s    *irb_ire;       /* First ire in this bucket */
1268                                         /* Should be first in this struct */
1269         krwlock_t       irb_lock;       /* Protect this bucket */
1270         uint_t          irb_refcnt;     /* Protected by irb_lock */
1271         uchar_t         irb_marks;      /* CONDEMNED ires in this bucket ? */
1272 #define IRB_MARK_CONDEMNED      0x0001  /* Contains some IRE_IS_CONDEMNED */
1273 #define IRB_MARK_DYNAMIC        0x0002  /* Dynamically allocated */
1274         /* Once IPv6 uses radix then IRB_MARK_DYNAMIC will be always be set */
1275         uint_t          irb_ire_cnt;    /* Num of active IRE in this bucket */
1276         int             irb_nire;       /* Num of ftable ire's that ref irb */
1277         ip_stack_t      *irb_ipst;      /* Does not have a netstack_hold */
1278 } irb_t;

1280 /*
1281  * This is the structure used to store the multicast physical addresses
1282  * that an interface has joined.
1283  * The refcnt keeps track of the number of multicast IP addresses mapping
1284  * to a physical multicast address.
1285  */
1286 typedef struct multiphysaddr_s {
1287         struct  multiphysaddr_s *mpa_next;
1288         char    mpa_addr[IP_MAX_HW_LEN];
1289         int     mpa_refcnt;
1290 } multiphysaddr_t;

1292 #define IRB2RT(irb)     (rt_t *)((caddr_t)(irb) - offsetof(rt_t, rt_irb))

1294 /* Forward declarations */
1295 struct dce_s;
1296 typedef struct dce_s dce_t;
1297 struct ire_s;
1298 typedef struct ire_s ire_t;
1299 struct ncec_s;
1300 typedef struct ncec_s ncec_t;
1301 struct nce_s;
1302 typedef struct nce_s nce_t;
1303 struct ip_recv_attr_s;
1304 typedef struct ip_recv_attr_s ip_recv_attr_t;
1305 struct ip_xmit_attr_s;
1306 typedef struct ip_xmit_attr_s ip_xmit_attr_t;

1308 struct tsol_ire_gw_secattr_s;
1309 typedef struct tsol_ire_gw_secattr_s tsol_ire_gw_secattr_t;

1311 /*
1312  * This is a structure for a one-element route cache that is passed
1313  * by reference between ip_input and ill_inputfn.
1314  */
1315 typedef struct {
```

```
1316            ire_t              *rtc_ire;
1317            ipaddr_t           rtc_ipaddr;
1318            in6_addr_t         rtc_ip6addr;
1319 } rtc_t;

1321 /*
1322  * Note: Temporarily use 64 bits, and will probably go back to 32 bits after
1323  * more cleanup work is done.
1324  */
1325 typedef uint64_t iaflags_t;

1327 /* The ill input function pointer type */
1328 typedef void (*pfillinput_t)(mblk_t *, void *, void *, ip_recv_attr_t *,
1329     rtc_t *);

1331 /* The ire receive function pointer type */
1332 typedef void (*pfirerecv_t)(ire_t *, mblk_t *, void *, ip_recv_attr_t *);

1334 /* The ire send and postfrag function pointer types */
1335 typedef int (*pfiresend_t)(ire_t *, mblk_t *, void *,
1336     ip_xmit_attr_t *, uint32_t *);
1337 typedef int (*pfirepostfrag_t)(mblk_t *, nce_t *, iaflags_t, uint_t, uint32_t,
1338     zoneid_t, zoneid_t, uintptr_t *);


1341 #define IP_V4_G_HEAD    0
1342 #define IP_V6_G_HEAD    1

1344 #define MAX_G_HEADS     2

1346 /*
1347  * unpadded ill_if structure
1348  */
1349 struct  _ill_if_s_ {
1350         union ill_if_u *illif_next;
1351         union ill_if_u *illif_prev;
1352         avl_tree_t      illif_avl_by_ppa;       /* AVL tree sorted on ppa */
1353         vmem_t         *illif_ppa_arena;        /* ppa index space */
1354         uint16_t        illif_mcast_v1;         /* hints for          */
1355         uint16_t        illif_mcast_v2;         /* [igmp|mld]_slowtimo    */
1356         int             illif_name_len;         /* name length */
1357         char            illif_name[LIFNAMSIZ];  /* name of interface type */
1358 };

1360 /* cache aligned ill_if structure */
1361 typedef union   ill_if_u {
1362         struct  _ill_if_s_ ill_if_s;
1363         char    illif_filler[CACHE_ALIGN(_ill_if_s_)];
1364 } ill_if_t;

1366 #define illif_next              ill_if_s.illif_next
1367 #define illif_prev              ill_if_s.illif_prev
1368 #define illif_avl_by_ppa        ill_if_s.illif_avl_by_ppa
1369 #define illif_ppa_arena         ill_if_s.illif_ppa_arena
1370 #define illif_mcast_v1          ill_if_s.illif_mcast_v1
1371 #define illif_mcast_v2          ill_if_s.illif_mcast_v2
1372 #define illif_name              ill_if_s.illif_name
1373 #define illif_name_len          ill_if_s.illif_name_len

1375 typedef struct ill_walk_context_s {
1376         int     ctx_current_list; /* current list being searched */
1377         int     ctx_last_list;  /* last list to search */
1378 } ill_walk_context_t;

1380 /*
1381  * ill_g_heads structure, one for IPV4 and one for IPV6
```

```
1382  */
1383 struct _ill_g_head_s_ {
1384         ill_if_t        *ill_g_list_head;
1385         ill_if_t        *ill_g_list_tail;
1386 };

1388 typedef union ill_g_head_u {
1389         struct _ill_g_head_s_ ill_g_head_s;
1390         char    ill_g_head_filler[CACHE_ALIGN(_ill_g_head_s_)];
1391 } ill_g_head_t;

1393 #define ill_g_list_head ill_g_head_s.ill_g_list_head
1394 #define ill_g_list_tail ill_g_head_s.ill_g_list_tail

1396 #define IP_V4_ILL_G_LIST(ipst)  \
1397         (ipst)->ips_ill_g_heads[IP_V4_G_HEAD].ill_g_list_head
1398 #define IP_V6_ILL_G_LIST(ipst)  \
1399         (ipst)->ips_ill_g_heads[IP_V6_G_HEAD].ill_g_list_head
1400 #define IP_VX_ILL_G_LIST(i, ipst)       \
1401         (ipst)->ips_ill_g_heads[i].ill_g_list_head

1403 #define ILL_START_WALK_V4(ctx_ptr, ipst)        \
1404         ill_first(IP_V4_G_HEAD, IP_V4_G_HEAD, ctx_ptr, ipst)
1405 #define ILL_START_WALK_V6(ctx_ptr, ipst)        \
1406         ill_first(IP_V6_G_HEAD, IP_V6_G_HEAD, ctx_ptr, ipst)
1407 #define ILL_START_WALK_ALL(ctx_ptr, ipst)       \
1408         ill_first(MAX_G_HEADS, MAX_G_HEADS, ctx_ptr, ipst)

1410 /*
1411  * Capabilities, possible flags for ill_capabilities.
1412  */
1413 #define ILL_CAPAB_LSO           0x04            /* Large Send Offload */
1414 #define ILL_CAPAB_HCKSUM        0x08            /* Hardware checksumming */
1415 #define ILL_CAPAB_ZEROCOPY      0x10            /* Zero-copy */
1416 #define ILL_CAPAB_DLD           0x20            /* DLD capabilities */
1417 #define ILL_CAPAB_DLD_POLL      0x40            /* Polling */
1418 #define ILL_CAPAB_DLD_DIRECT    0x80            /* Direct function call */

1420 /*
1421  * Per-ill Hardware Checksumming capbilities.
1422  */
1423 typedef struct ill_hcksum_capab_s ill_hcksum_capab_t;

1425 /*
1426  * Per-ill Zero-copy capabilities.
1427  */
1428 typedef struct ill_zerocopy_capab_s ill_zerocopy_capab_t;

1430 /*
1431  * DLD capbilities.
1432  */
1433 typedef struct ill_dld_capab_s ill_dld_capab_t;

1435 /*
1436  * Per-ill polling resource map.
1437  */
1438 typedef struct ill_rx_ring ill_rx_ring_t;

1440 /*
1441  * Per-ill Large Send Offload capabilities.
1442  */
1443 typedef struct ill_lso_capab_s ill_lso_capab_t;

1445 /* The following are ill_state_flags */
1446 #define ILL_LL_SUBNET_PENDING   0x01    /* Waiting for DL_INFO_ACK from drv */
1447 #define ILL_CONDEMNED           0x02    /* No more new ref's to the ILL */
```

```
1448 #define ILL_DL_UNBIND_IN_PROGRESS       0x04    /* UNBIND_REQ is sent */
1449 /*
1450  * ILL_DOWN_IN_PROGRESS is set to ensure the following:
1451  * - no packets are sent to the driver after the DL_UNBIND_REQ is sent,
1452  * - no longstanding references will be acquired on objects that are being
1453  *   brought down.
1454  */
1455 #define ILL_DOWN_IN_PROGRESS    0x08

1457 /* Is this an ILL whose source address is used by other ILL's ? */
1458 #define IS_USESRC_ILL(ill)                              \
1459         (((ill)->ill_usesrc_ifindex == 0) &&    \
1460         ((ill)->ill_usesrc_grp_next != NULL))

1462 /* Is this a client/consumer of the usesrc ILL ? */
1463 #define IS_USESRC_CLI_ILL(ill)                          \
1464         (((ill)->ill_usesrc_ifindex != 0) &&    \
1465         ((ill)->ill_usesrc_grp_next != NULL))

1467 /* Is this an virtual network interface (vni) ILL ? */
1468 #define IS_VNI(ill)                                                     \
1469         (((ill)->ill_phyint->phyint_flags & (PHYI_LOOPBACK|PHYI_VIRTUAL)) == \
1470         PHYI_VIRTUAL)

1472 /* Is this a loopback ILL? */
1473 #define IS_LOOPBACK(ill) \
1474         ((ill)->ill_phyint->phyint_flags & PHYI_LOOPBACK)

1476 /* Is this an IPMP meta-interface ILL? */
1477 #define IS_IPMP(ill)                                                    \
1478         ((ill)->ill_phyint->phyint_flags & PHYI_IPMP)

1480 /* Is this ILL under an IPMP meta-interface? (aka "in a group?") */
1481 #define IS_UNDER_IPMP(ill)                                              \
1482         ((ill)->ill_grp != NULL && !IS_IPMP(ill))

1484 /* Is ill1 in the same illgrp as ill2? */
1485 #define IS_IN_SAME_ILLGRP(ill1, ill2)                                   \
1486         ((ill1)->ill_grp != NULL && ((ill1)->ill_grp == (ill2)->ill_grp))

1488 /* Is ill1 on the same LAN as ill2? */
1489 #define IS_ON_SAME_LAN(ill1, ill2)                                      \
1490         ((ill1) == (ill2) || IS_IN_SAME_ILLGRP(ill1, ill2))

1492 #define ILL_OTHER(ill)                                                  \
1493         ((ill)->ill_isv6 ? (ill)->ill_phyint->phyint_illv4 :            \
1494             (ill)->ill_phyint->phyint_illv6)
1496 /*
1497  * IPMP group ILL state structure -- up to two per IPMP group (V4 and V6).
1498  * Created when the V4 and/or V6 IPMP meta-interface is I_PLINK'd.  It is
1499  * guaranteed to persist while there are interfaces in the group.
1500  * In general, most fields are accessed outside of the IPSQ (e.g., in the
1501  * datapath), and thus use locks in addition to the IPSQ for protection.
1502  *
1503  * synchronization:             read                    write
1504  *
1505  *      ig_if                   ipsq or ill_g_lock      ipsq and ill_g_lock
1506  *      ig_actif                ipsq or ipmp_lock       ipsq and ipmp_lock
1507  *      ig_nactif               ipsq or ipmp_lock       ipsq and ipmp_lock
1508  *      ig_next_ill             ipsq or ipmp_lock       ipsq and ipmp_lock
1509  *      ig_ipmp_ill             write once              write once
1510  *      ig_cast_ill             ipsq or ipmp_lock       ipsq and ipmp_lock
1511  *      ig_arpent               ipsq                    ipsq
1512  *      ig_mtu                  ipsq                    ipsq
1513  *      ig_mc_mtu               ipsq                    ipsq
```

```
1514  */
1515 typedef struct ipmp_illgrp_s {
1516         list_t          ig_if;          /* list of all interfaces */
1517         list_t          ig_actif;       /* list of active interfaces */
1518         uint_t          ig_nactif;      /* number of active interfaces */
1519         struct ill_s    *ig_next_ill;   /* next active interface to use */
1520         struct ill_s    *ig_ipmp_ill;   /* backpointer to IPMP meta-interface */
1521         struct ill_s    *ig_cast_ill;   /* nominated ill for multi/broadcast */
1522         list_t          ig_arpent;      /* list of ARP entries */
1523         uint_t          ig_mtu;         /* ig_ipmp_ill->ill_mtu */
1524         uint_t          ig_mc_mtu;      /* ig_ipmp_ill->ill_mc_mtu */
1525 } ipmp_illgrp_t;

1527 /*
1528  * IPMP group state structure -- one per IPMP group.  Created when the
1529  * IPMP meta-interface is plumbed; it is guaranteed to persist while there
1530  * are interfaces in it.
1531  *
1532  * ipmp_grp_t synchronization:          read                    write
1533  *
1534  *      gr_name                         ipmp_lock               ipmp_lock
1535  *      gr_ifname                       write once              write once
1536  *      gr_mactype                      ipmp_lock               ipmp_lock
1537  *      gr_phyint                       write once              write once
1538  *      gr_nif                          ipmp_lock               ipmp_lock
1539  *      gr_nactif                       ipsq                    ipsq
1540  *      gr_v4                           ipmp_lock               ipmp_lock
1541  *      gr_v6                           ipmp_lock               ipmp_lock
1542  *      gr_nv4                          ipmp_lock               ipmp_lock
1543  *      gr_nv6                          ipmp_lock               ipmp_lock
1544  *      gr_pendv4                       ipmp_lock               ipmp_lock
1545  *      gr_pendv6                       ipmp_lock               ipmp_lock
1546  *      gr_linkdownmp                   ipsq                    ipsq
1547  *      gr_ksp                          ipmp_lock               ipmp_lock
1548  *      gr_kstats0                      atomic                  atomic
1549  */
1550 typedef struct ipmp_grp_s {
1551         char            gr_name[LIFGRNAMSIZ];   /* group name */
1552         char            gr_ifname[LIFNAMSIZ];   /* interface name */
1553         t_uscalar_t     gr_mactype;     /* DLPI mactype of group */
1554         phyint_t        *gr_phyint;     /* IPMP group phyint */
1555         uint_t          gr_nif;         /* number of interfaces in group */
1556         uint_t          gr_nactif;      /* number of active interfaces */
1557         ipmp_illgrp_t   *gr_v4;         /* V4 group information */
1558         ipmp_illgrp_t   *gr_v6;         /* V6 group information */
1559         uint_t          gr_nv4;         /* number of ills in V4 group */
1560         uint_t          gr_nv6;         /* number of ills in V6 group */
1561         uint_t          gr_pendv4;      /* number of pending ills in V4 group */
1562         uint_t          gr_pendv6;      /* number of pending ills in V6 group */
1563         mblk_t          *gr_linkdownmp; /* message used to bring link down */
1564         kstat_t         *gr_ksp;        /* group kstat pointer */
1565         uint64_t        gr_kstats0[IPMP_KSTAT_MAX]; /* baseline group kstats */
1566 } ipmp_grp_t;

1568 /*
1569  * IPMP ARP entry -- one per SIOCS*ARP entry tied to the group.  Used to keep
1570  * ARP up-to-date as the active set of interfaces in the group changes.
1571  */
1572 typedef struct ipmp_arpent_s {
1573         ipaddr_t        ia_ipaddr;      /* IP address for this entry */
1574         boolean_t       ia_proxyarp;    /* proxy ARP entry? */
1575         boolean_t       ia_notified;    /* ARP notified about this entry? */
1576         list_node_t     ia_node;        /* next ARP entry in list */
1577         uint16_t        ia_flags;       /* nce_flags for the address */
1578         size_t          ia_lladdr_len;
1579         uchar_t         *ia_lladdr;
```

```
1580 } ipmp_arpent_t;

1582 struct arl_s;

1584 /*
1585  * Per-ill capabilities.
1586  */
1587 struct ill_hcksum_capab_s {
1588         uint_t  ill_hcksum_version;     /* interface version */
1589         uint_t  ill_hcksum_txflags;     /* capabilities on transmit */
1590 };

1592 struct ill_zerocopy_capab_s {
1593         uint_t  ill_zerocopy_version;   /* interface version */
1594         uint_t  ill_zerocopy_flags;     /* capabilities */
1595 };

1597 struct ill_lso_capab_s {
1598         uint_t  ill_lso_flags;          /* capabilities */
1599         uint_t  ill_lso_max;            /* maximum size of payload */
1600 };
1602 /*
1603  * IP Lower level Structure.
1604  * Instance data structure in ip_open when there is a device below us.
1605  */
1606 typedef struct ill_s {
1607         pfillinput_t ill_inputfn;       /* Fast input function selector */
1608         ill_if_t *ill_ifptr;            /* pointer to interface type */
1609         queue_t *ill_rq;                /* Read queue. */
1610         queue_t *ill_wq;                /* Write queue. */

1612         int     ill_error;              /* Error value sent up by device. */

1614         ipif_t  *ill_ipif;              /* Interface chain for this ILL. */

1616         uint_t  ill_ipif_up_count;      /* Number of IPIFs currently up. */
1617         uint_t  ill_max_frag;           /* Max IDU from DLPI. */
1618         uint_t  ill_current_frag;       /* Current IDU from DLPI. */
1619         uint_t  ill_mtu;                /* User-specified MTU; SIOCSLIFMTU */
1620         uint_t  ill_mc_mtu;             /* MTU for multi/broadcast */
1621         uint_t  ill_metric;             /* BSD if metric, for compatibility. */
1622         char    *ill_name;              /* Our name. */
1623         uint_t  ill_ipif_dup_count;     /* Number of duplicate addresses. */
1624         uint_t  ill_name_length;        /* Name length, incl. terminator. */
1625         uint_t  ill_net_type;           /* IRE_IF_RESOLVER/IRE_IF_NORESOLVER. */
1626         /*
1627          * Physical Point of Attachment num.  If DLPI style 1 provider
1628          * then this is derived from the devname.
1629          */
1630         uint_t  ill_ppa;
1631         t_uscalar_t     ill_sap;
1632         t_scalar_t      ill_sap_length; /* Including sign (for position) */
1633         uint_t  ill_phys_addr_length;   /* Excluding the sap. */
1634         uint_t  ill_bcast_addr_length;  /* Only set when the DL provider */
1635                                         /* supports broadcast. */
1636         t_uscalar_t     ill_mactype;
1637         uint8_t *ill_frag_ptr;          /* Reassembly state. */
1638         timeout_id_t ill_frag_timer_id; /* timeout id for the frag timer */
1639         ipfb_t  *ill_frag_hash_tbl;     /* Fragment hash list head. */

1641         krwlock_t ill_mcast_lock;       /* Protects multicast state */
1642         kmutex_t ill_mcast_serializer;  /* Serialize across ilg and ilm state */
1643         ilm_t   *ill_ilm;               /* Multicast membership for ill */
1644         uint_t  ill_global_timer;       /* for IGMPv3/MLDv2 general queries */
1645         int     ill_mcast_type;         /* type of router which is querier */
```

```
1646                                         /* on this interface */
1647         uint16_t ill_mcast_v1_time;     /* # slow timeouts since last v1 qry */
1648         uint16_t ill_mcast_v2_time;     /* # slow timeouts since last v2 qry */
1649         uint8_t ill_mcast_v1_tset;      /* 1 => timer is set; 0 => not set */
1650         uint8_t ill_mcast_v2_tset;      /* 1 => timer is set; 0 => not set */

1652         uint8_t ill_mcast_rv;           /* IGMPv3/MLDv2 robustness variable */
1653         int     ill_mcast_qi;           /* IGMPv3/MLDv2 query interval var */

1655         /*
1656          * All non-NULL cells between 'ill_first_mp_to_free' and
1657          * 'ill_last_mp_to_free' are freed in ill_delete.
1658          */
1659 #define ill_first_mp_to_free    ill_bcast_mp
1660         mblk_t  *ill_bcast_mp;          /* DLPI header for broadcasts. */
1661         mblk_t  *ill_unbind_mp;         /* unbind mp from ill_dl_up() */
1662         mblk_t  *ill_promiscoff_mp;     /* for ill_leave_allmulti() */
1663         mblk_t  *ill_dlpi_deferred;     /* b_next chain of control messages */
1664         mblk_t  *ill_dest_addr_mp;      /* mblk which holds ill_dest_addr */
1665         mblk_t  *ill_replumb_mp;        /* replumb mp from ill_replumb() */
1666         mblk_t  *ill_phys_addr_mp;      /* mblk which holds ill_phys_addr */
1667         mblk_t  *ill_mcast_deferred;    /* b_next chain of IGMP/MLD packets */
1668 #define ill_last_mp_to_free     ill_mcast_deferred

1670         cred_t  *ill_credp;             /* opener's credentials */
1671         uint8_t *ill_phys_addr;         /* ill_phys_addr_mp->b_rptr + off */
1672         uint8_t *ill_dest_addr;         /* ill_dest_addr_mp->b_rptr + off */

1674         uint_t  ill_state_flags;        /* see ILL_* flags above */

1676         /* Following bit fields protected by ipsq_t */
1677         uint_t
1678                 ill_needs_attach : 1,
1679                 ill_reserved : 1,
1680                 ill_isv6 : 1,
1681                 ill_dlpi_style_set : 1,

1683                 ill_ifname_pending : 1,
1684                 ill_logical_down : 1,
1685                 ill_dl_up : 1,
1686                 ill_up_ipifs : 1,

1688                 ill_note_link : 1,      /* supports link-up notification */
1689                 ill_capab_reneg : 1, /* capability renegotiation to be done */
1690                 ill_dld_capab_inprog : 1, /* direct dld capab call in prog */
1691                 ill_need_recover_multicast : 1,

1693                 ill_replumbing : 1,
1694                 ill_arl_dlpi_pending : 1,
1695                 ill_grp_pending : 1,

1697                 ill_pad_to_bit_31 : 17;

1699         /* Following bit fields protected by ill_lock */
1700         uint_t
1701                 ill_fragtimer_executing : 1,
1702                 ill_fragtimer_needrestart : 1,
1703                 ill_manual_token : 1,   /* system won't override ill_token */
1704                 /*
1705                  * ill_manual_linklocal : system will not change the
1706                  * linklocal whenever ill_token changes.
1707                  */
1708                 ill_manual_linklocal : 1,

1710                 ill_manual_dst_linklocal : 1, /* same for pt-pt dst linklocal */
```

```
1712                     ill_pad_bit_31 : 27;

1714         /*
1715          * Used in SIOCSIFMUXID and SIOCGIFMUXID for 'ifconfig unplumb'.
1716          */
1717         int     ill_muxid;              /* muxid returned from plink */

1719         /* Used for IP frag reassembly throttling on a per ILL basis.  */
1720         uint_t  ill_ipf_gen;            /* Generation of next fragment queue */
1721         uint_t  ill_frag_count;         /* Count of all reassembly mblk bytes */
1722         uint_t  ill_frag_free_num_pkts;  /* num of fragmented packets to free */
1723         clock_t ill_last_frag_clean_time; /* time when frag's were pruned */
1724         int     ill_type;               /* From <net/if_types.h> */
1725         uint_t  ill_dlpi_multicast_state;       /* See below IDS_* */
1726         uint_t  ill_dlpi_fastpath_state;        /* See below IDS_* */

1728         /*
1729          * Capabilities related fields.
1730          */
1731         uint_t  ill_dlpi_capab_state;   /* State of capability query, IDCS_* */
1732         uint_t  ill_capab_pending_cnt;
1733         uint64_t ill_capabilities;      /* Enabled capabilities, ILL_CAPAB_* */
1734         ill_hcksum_capab_t *ill_hcksum_capab; /* H/W cksumming capabilities */
1735         ill_zerocopy_capab_t *ill_zerocopy_capab; /* Zero-copy capabilities */
1736         ill_dld_capab_t *ill_dld_capab; /* DLD capabilities */
1737         ill_lso_capab_t *ill_lso_capab; /* Large Segment Offload capabilities */
1738         mblk_t  *ill_capab_reset_mp;    /* Preallocated mblk for capab reset */

1740         uint8_t ill_max_hops;   /* Maximum hops for any logical interface */
1741         uint_t  ill_user_mtu;   /* User-specified MTU via SIOCSLIFLNKINFO */
1742         uint32_t ill_reachable_time;    /* Value for ND algorithm in msec */
1743         uint32_t ill_reachable_retrans_time; /* Value for ND algorithm msec */
1744         uint_t  ill_max_buf;            /* Max # of req to buffer for ND */
1745         in6_addr_t      ill_token;      /* IPv6 interface id */
1746         in6_addr_t      ill_dest_token; /* Destination IPv6 interface id */
1747         uint_t          ill_token_length;
1748         uint32_t        ill_xmit_count;         /* ndp max multicast xmits */
1749         mib2_ipIfStatsEntry_t   *ill_ip_mib;    /* ver indep. interface mib */
1750         mib2_ipv6IfIcmpEntry_t  *ill_icmp6_mib; /* Per interface mib */

1752         phyint_t                *ill_phyint;
1753         uint64_t                ill_flags;

1755         kmutex_t        ill_lock;       /* Please see table below */
1756         /*
1757          * The ill_nd_lla* fields handle the link layer address option
1758          * from neighbor discovery. This is used for external IPv6
1759          * address resolution.
1760          */
1761         mblk_t          *ill_nd_lla_mp; /* mblk which holds ill_nd_lla */
1762         uint8_t         *ill_nd_lla;    /* Link Layer Address */
1763         uint_t          ill_nd_lla_len; /* Link Layer Address length */
1764         /*
1765          * We have 4 phys_addr_req's sent down. This field keeps track
1766          * of which one is pending.
1767          */
1768         t_uscalar_t     ill_phys_addr_pend; /* which dl_phys_addr_req pending */
1769         /*
1770          * Used to save errors that occur during plumbing
1771          */
1772         uint_t          ill_ifname_pending_err;
1773         avl_node_t      ill_avl_byppa; /* avl node based on ppa */
1774         list_t          ill_nce; /* pointer to nce_s list */
1775         uint_t          ill_refcnt;     /* active refcnt by threads */
1776         uint_t          ill_ire_cnt;    /* ires associated with this ill */
1777         kcondvar_t      ill_cv;
```

```
1778         uint_t          ill_ncec_cnt;   /* ncecs associated with this ill */
1779         uint_t          ill_nce_cnt;    /* nces associated with this ill */
1780         uint_t          ill_waiters;    /* threads waiting in ipsq_enter */
1781         /*
1782          * Contains the upper read queue pointer of the module immediately
1783          * beneath IP.  This field allows IP to validate sub-capability
1784          * acknowledgments coming up from downstream.
1785          */
1786         queue_t         *ill_lmod_rq;   /* read queue pointer of module below */
1787         uint_t          ill_lmod_cnt;   /* number of modules beneath IP */
1788         ip_m_t          *ill_media;     /* media specific params/functions */
1789         t_uscalar_t     ill_dlpi_pending; /* Last DLPI primitive issued */
1790         uint_t          ill_usesrc_ifindex; /* use src addr from this ILL */
1791         struct ill_s    *ill_usesrc_grp_next; /* Next ILL in the usesrc group */
1792         boolean_t       ill_trace_disable;      /* True when alloc fails */
1793         zoneid_t        ill_zoneid;
1794         ip_stack_t      *ill_ipst;      /* Corresponds to a netstack_hold */
1795         uint32_t        ill_dhcpinit;   /* IP_DHCPINIT_IFs for ill */
1796         void            *ill_flownotify_mh; /* Tx flow ctl, mac cb handle */
1797         uint_t          ill_ilm_cnt;    /* ilms referencing this ill */
1798         uint_t          ill_ipallmulti_cnt; /* ip_join_allmulti() calls */
1799         ilm_t           *ill_ipallmulti_ilm;

1801         mblk_t          *ill_saved_ire_mp; /* Allocated for each extra IRE */
1802                                         /* with ire_ill set so they can */
1803                                         /* survive the ill going down and up. */
1804         kmutex_t        ill_saved_ire_lock; /* Protects ill_saved_ire_mp, cnt */
1805         uint_t          ill_saved_ire_cnt;      /* # entries */
1806         struct arl_ill_common_s    *ill_common;
1807         ire_t           *ill_ire_multicast; /* IRE_MULTICAST for ill */
1808         clock_t         ill_defend_start;   /* start of 1 hour period */
1809         uint_t          ill_defend_count;   /* # of announce/defends per ill */
1810         /*
1811          * IPMP fields.
1812          */
1813         ipmp_illgrp_t   *ill_grp;       /* IPMP group information */
1814         list_node_t     ill_actnode;    /* next active ill in group */
1815         list_node_t     ill_grpnode;    /* next ill in group */
1816         ipif_t          *ill_src_ipif;  /* source address selection rotor */
1817         ipif_t          *ill_move_ipif; /* ipif awaiting move to new ill */
1818         boolean_t       ill_nom_cast;   /* nominated for mcast/bcast */
1819         uint_t          ill_bound_cnt;  /* # of data addresses bound to ill */
1820         ipif_t          *ill_bound_ipif; /* ipif chain bound to ill */
1821         timeout_id_t    ill_refresh_tid; /* ill refresh retry timeout id */

1823         uint32_t        ill_mrouter_cnt; /* mrouter allmulti joins */
1824         uint32_t        ill_allowed_ips_cnt;
1825         in6_addr_t      *ill_allowed_ips;

1827         /* list of multicast physical addresses joined on this ill */
1828         multiphysaddr_t *ill_mphysaddr_list;
1829 } ill_t;

1831 /*
1832  * ILL_FREE_OK() means that there are no incoming pointer references
1833  * to the ill.
1834  */
1835 #define ILL_FREE_OK(ill)                                                \
1836         ((ill)->ill_ire_cnt == 0 && (ill)->ill_ilm_cnt == 0 &&  \
1837         (ill)->ill_ncec_cnt == 0 && (ill)->ill_nce_cnt == 0)

1839 /*
1840  * An ipif/ill can be marked down only when the ire and ncec references
1841  * to that ipif/ill goes to zero. ILL_DOWN_OK() is a necessary condition
1842  * quiescence checks. See comments above IPIF_DOWN_OK for details
1843  * on why ires and nces are selectively considered for this macro.
```

```
1844  */
1845  #define ILL_DOWN_OK(ill)                                              \
1846          (ill->ill_ire_cnt == 0 && ill->ill_ncec_cnt == 0 &&     \
1847          ill->ill_nce_cnt == 0)
1848
1849  /*
1850   * The following table lists the protection levels of the various members
1851   * of the ill_t. Same notation as that used for ipif_t above is used.
1852   *
1853   *                              Write                   Read
1854   *
1855   * ill_ifptr                    ill_g_lock + s          Write once
1856   * ill_rq                       ipsq                    Write once
1857   * ill_wq                       ipsq                    Write once
1858   *
1859   * ill_error                    ipsq                    None
1860   * ill_ipif                     ill_g_lock + ipsq       ill_g_lock OR ipsq
1861   * ill_ipif_up_count            ill_lock + ipsq         ill_lock OR ipsq
1862   * ill_max_frag                 ill_lock                ill_lock
1863   * ill_current_frag             ill_lock                ill_lock
1864   *
1865   * ill_name                     ill_g_lock + ipsq       Write once
1866   * ill_name_length              ill_g_lock + ipsq       Write once
1867   * ill_ndd_name                 ipsq                    Write once
1868   * ill_net_type                 ipsq                    Write once
1869   * ill_ppa                      ill_g_lock + ipsq       Write once
1870   * ill_sap                      ipsq + down ill         Write once
1871   * ill_sap_length               ipsq + down ill         Write once
1872   * ill_phys_addr_length         ipsq + down ill         Write once
1873   *
1874   * ill_bcast_addr_length        ipsq                    ipsq
1875   * ill_mactype                  ipsq                    ipsq
1876   * ill_frag_ptr                 ipsq                    ipsq
1877   *
1878   * ill_frag_timer_id            ill_lock                ill_lock
1879   * ill_frag_hash_tbl            ipsq                    up ill
1880   * ill_ilm                      ill_mcast_lock(WRITER)  ill_mcast_lock(READER)
1881   * ill_global_timer             ill_mcast_lock(WRITER)  ill_mcast_lock(READER)
1882   * ill_mcast_type               ill_mcast_lock(WRITER)  ill_mcast_lock(READER)
1883   * ill_mcast_v1_time            ill_mcast_lock(WRITER)  ill_mcast_lock(READER)
1884   * ill_mcast_v2_time            ill_mcast_lock(WRITER)  ill_mcast_lock(READER)
1885   * ill_mcast_v1_tset            ill_mcast_lock(WRITER)  ill_mcast_lock(READER)
1886   * ill_mcast_v2_tset            ill_mcast_lock(WRITER)  ill_mcast_lock(READER)
1887   * ill_mcast_rv                 ill_mcast_lock(WRITER)  ill_mcast_lock(READER)
1888   * ill_mcast_qi                 ill_mcast_lock(WRITER)  ill_mcast_lock(READER)
1889   *
1890   * ill_down_mp                  ipsq                    ipsq
1891   * ill_dlpi_deferred            ill_lock                ill_lock
1892   * ill_dlpi_pending             ipsq + ill_lock         ipsq or ill_lock or
1893   *                                                      absence of ipsq writer.
1894   * ill_phys_addr_mp             ipsq + down ill         only when ill is up
1895   * ill_mcast_deferred           ill_lock                ill_lock
1896   * ill_phys_addr                ipsq + down ill         only when ill is up
1897   * ill_dest_addr_mp             ipsq + down ill         only when ill is up
1898   * ill_dest_addr                ipsq + down ill         only when ill is up
1899   *
1900   * ill_state_flags              ill_lock                ill_lock
1901   * exclusive bit flags          ipsq_t                  ipsq_t
1902   * shared bit flags             ill_lock                ill_lock
1903   *
1904   * ill_muxid                    ipsq                    Not atomic
1905   *
1906   * ill_ipf_gen                  Not atomic
1907   * ill_frag_count               atomics                 atomics
1908   * ill_type                     ipsq + down ill         only when ill is up
1909   * ill_dlpi_multicast_state     ill_lock                ill_lock
```

```
1910   * ill_dlpi_fastpath_state      ill_lock                ill_lock
1911   * ill_dlpi_capab_state         ipsq                    ipsq
1912   * ill_max_hops                 ipsq                    Not atomic
1913   *
1914   * ill_mtu                      ill_lock                None
1915   * ill_mc_mtu                   ill_lock                None
1916   *
1917   * ill_user_mtu                 ipsq + ill_lock         ill_lock
1918   * ill_reachable_time           ipsq + ill_lock         ill_lock
1919   * ill_reachable_retrans_time   ipsq + ill_lock         ill_lock
1920   * ill_max_buf                  ipsq + ill_lock         ill_lock
1921   *
1922   * Next 2 fields need ill_lock because of the get ioctls. They should not
1923   * report partially updated results without executing in the ipsq.
1924   * ill_token                    ipsq + ill_lock         ill_lock
1925   * ill_token_length             ipsq + ill_lock         ill_lock
1926   * ill_dest_token               ipsq + down ill         only when ill is up
1927   * ill_xmit_count               ipsq + down ill         write once
1928   * ill_ip6_mib                  ipsq + down ill         only when ill is up
1929   * ill_icmp6_mib                ipsq + down ill         only when ill is up
1930   *
1931   * ill_phyint                   ipsq, ill_g_lock, ill_lock      Any of them
1932   * ill_flags                    ill_lock                ill_lock
1933   * ill_nd_lla_mp                ipsq + down ill         only when ill is up
1934   * ill_nd_lla                   ipsq + down ill         only when ill is up
1935   * ill_nd_lla_len               ipsq + down ill         only when ill is up
1936   * ill_phys_addr_pend           ipsq + down ill         only when ill is up
1937   * ill_ifname_pending_err       ipsq                    ipsq
1938   * ill_avl_byppa                ipsq, ill_g_lock        write once
1939   *
1940   * ill_fastpath_list            ill_lock                ill_lock
1941   * ill_refcnt                   ill_lock                ill_lock
1942   * ill_ire_cnt                  ill_lock                ill_lock
1943   * ill_cv                       ill_lock                ill_lock
1944   * ill_ncec_cnt                 ill_lock                ill_lock
1945   * ill_nce_cnt                  ill_lock                ill_lock
1946   * ill_ilm_cnt                  ill_lock                ill_lock
1947   * ill_src_ipif                 ill_g_lock              ill_g_lock
1948   * ill_trace                    ill_lock                ill_lock
1949   * ill_usesrc_grp_next          ill_g_usesrc_lock       ill_g_usesrc_lock
1950   * ill_dhcpinit                 atomics                 atomics
1951   * ill_flownotify_mh            write once              write once
1952   * ill_capab_pending_cnt        ipsq                    ipsq
1953   * ill_ipallmulti_cnt           ill_lock                ill_lock
1954   * ill_ipallmulti_ilm           ill_lock                ill_lock
1955   * ill_saved_ire_mp             ill_saved_ire_lock      ill_saved_ire_lock
1956   * ill_saved_ire_cnt            ill_saved_ire_lock      ill_saved_ire_lock
1957   * ill_arl                      ???                     ???
1958   * ill_ire_multicast            ipsq + quiescent        none
1959   * ill_bound_ipif               ipsq                    ipsq
1960   * ill_actnode                  ipsq + ipmp_lock        ipsq OR ipmp_lock
1961   * ill_grpnode                  ipsq + ill_g_lock       ipsq OR ill_g_lock
1962   * ill_src_ipif                 ill_g_lock              ill_g_lock
1963   * ill_move_ipif                ipsq                    ipsq
1964   * ill_nom_cast                 ipsq                    ipsq OR advisory
1965   * ill_refresh_tid              ill_lock                ill_lock
1966   * ill_grp (for IPMP ill)       write once              write once
1967   * ill_grp (for underlying ill) ipsq + ill_g_lock       ipsq OR ill_g_lock
1968   * ill_grp_pending              ill_mcast_serializer    ill_mcast_serializer
1969   * ill_mrouter_cnt              atomics                 atomics
1970   * ill_mphysaddr_list   ill_lock                ill_lock
1971   *
1972   * NOTE: It's OK to make heuristic decisions on an underlying interface
1973   *       by using IS_UNDER_IPMP() or comparing ill_grp's raw pointer value.
1974   */
```

```
1976 /*
1977  * For ioctl restart mechanism see ip_reprocess_ioctl()
1978  */
1979 struct ip_ioctl_cmd_s;

1981 typedef int (*ifunc_t)(ipif_t *, struct sockaddr_in *, queue_t *, mblk_t *,
1982     struct ip_ioctl_cmd_s *, void *);

1984 typedef struct ip_ioctl_cmd_s {
1985         int     ipi_cmd;
1986         size_t  ipi_copyin_size;
1987         uint_t  ipi_flags;
1988         uint_t  ipi_cmd_type;
1989         ifunc_t ipi_func;
1990         ifunc_t ipi_func_restart;
1991 } ip_ioctl_cmd_t;

1993 /*
1994  * ipi_cmd_type:
1995  *
1996  * IF_CMD               1       old style ifreq cmd
1997  * LIF_CMD              2       new style lifreq cmd
1998  * ARP_CMD              3       arpreq cmd
1999  * XARP_CMD             4       xarpreq cmd
2000  * MSFILT_CMD           5       multicast source filter cmd
2001  * MISC_CMD             6       misc cmd (not a more specific one above)
2002  */

2004 enum { IF_CMD = 1, LIF_CMD, ARP_CMD, XARP_CMD, MSFILT_CMD, MISC_CMD };

2006 #define IPI_DONTCARE   0       /* For ioctl encoded values that don't matter */

2008 /* Flag values in ipi_flags */
2009 #define IPI_PRIV       0x1     /* Root only command */
2010 #define IPI_MODOK      0x2     /* Permitted on mod instance of IP */
2011 #define IPI_WR         0x4     /* Need to grab writer access */
2012 #define IPI_GET_CMD    0x8     /* branch to mi_copyout on success */
2013 /*      unused         0x10    */
2014 #define IPI_NULL_BCONT 0x20    /* ioctl has not data and hence no b_cont */

2016 extern ip_ioctl_cmd_t   ip_ndx_ioctl_table[];
2017 extern ip_ioctl_cmd_t   ip_misc_ioctl_table[];
2018 extern int ip_ndx_ioctl_count;
2019 extern int ip_misc_ioctl_count;

2021 /* Passed down by ARP to IP during I_PLINK/I_PUNLINK */
2022 typedef struct ipmx_s {
2023         char    ipmx_name[LIFNAMSIZ];           /* if name */
2024         uint_t
2025                 ipmx_arpdev_stream : 1,         /* This is the arp stream */
2026                 ipmx_notused : 31;
2027 } ipmx_t;

2029 /*
2030  * State for detecting if a driver supports certain features.
2031  * Support for DL_ENABMULTI_REQ uses ill_dlpi_multicast_state.
2032  * Support for DLPI M_DATA fastpath uses ill_dlpi_fastpath_state.
2033  */
2034 #define IDS_UNKNOWN    0       /* No DLPI request sent */
2035 #define IDS_INPROGRESS 1       /* DLPI request sent */
2036 #define IDS_OK         2       /* DLPI request completed successfully */
2037 #define IDS_FAILED     3       /* DLPI request failed */

2039 /* Support for DL_CAPABILITY_REQ uses ill_dlpi_capab_state. */
2040 enum {
2041         IDCS_UNKNOWN,
```

```
2042         IDCS_PROBE_SENT,
2043         IDCS_OK,
2044         IDCS_RESET_SENT,
2045         IDCS_RENEG,
2046         IDCS_FAILED
2047 };

2049 /* Extended NDP Management Structure */
2050 typedef struct ipndp_s {
2051         ndgetf_t        ip_ndp_getf;
2052         ndsetf_t        ip_ndp_setf;
2053         caddr_t         ip_ndp_data;
2054         char            *ip_ndp_name;
2055 } ipndp_t;

2057 /* IXA Notification types */
2058 typedef enum {
2059         IXAN_LSO,       /* LSO capability change */
2060         IXAN_PMTU,      /* PMTU change */
2061         IXAN_ZCOPY      /* ZEROCOPY capability change */
2062 } ixa_notify_type_t;

2064 typedef uint_t ixa_notify_arg_t;

2066 typedef void    (*ixa_notify_t)(void *, ip_xmit_attr_t *ixa, ixa_notify_type_t,
2067     ixa_notify_arg_t);

2069 /*
2070  * Attribute flags that are common to the transmit and receive attributes
2071  */
2072 #define IAF_IS_IPV4             0x80000000      /* ipsec_*_v4 */
2073 #define IAF_TRUSTED_ICMP        0x40000000      /* ipsec_*_icmp_loopback */
2074 #define IAF_NO_LOOP_ZONEID_SET  0x20000000      /* Zone that shouldn't have */
2075                                                 /* a copy */
2076 #define IAF_LOOPBACK_COPY       0x10000000      /* For multi and broadcast */

2078 #define IAF_MASK                0xf0000000      /* Flags that are common */

2080 /*
2081  * Transmit side attributes used between the transport protocols and IP as
2082  * well as inside IP. It is also used to cache information in the conn_t i.e.
2083  * replaces conn_ire and the IPsec caching in the conn_t.
2084  */
2085 struct ip_xmit_attr_s {
2086         iaflags_t       ixa_flags;      /* IXAF_*. See below */

2088         uint32_t        ixa_free_flags; /* IXA_FREE_*. See below */
2089         uint32_t        ixa_refcnt;     /* Using atomics */

2091         /*
2092          * Always initialized independently of ixa_flags settings.
2093          * Used by ip_xmit so we keep them up front for cache locality.
2094          */
2095         uint32_t        ixa_xmit_hint;  /* For ECMP and GLD TX ring fanout */
2096         uint_t          ixa_pktlen;     /* Always set. For frag and stats */
2097         zoneid_t        ixa_zoneid;     /* Assumed always set */

2099         /* Always set for conn_ip_output(); might be stale */
2100         /*
2101          * Since TCP keeps the conn_t around past the process going away
2102          * we need to use the "notr" (e.g, ire_refhold_notr) for ixa_ire,
2103          * ixa_nce, and ixa_dce.
2104          */
2105         ire_t           *ixa_ire;       /* Forwarding table entry */
2106         uint_t          ixa_ire_generation;
2107         nce_t           *ixa_nce;       /* Neighbor cache entry */
```

```
2108         dce_t          *ixa_dce;       /* Destination cache entry */
2109         uint_t         ixa_dce_generation;
2110         uint_t         ixa_src_generation;     /* If IXAF_VERIFY_SOURCE */

2112         uint32_t       ixa_src_preferences;    /* prefs for src addr select */
2113         uint32_t       ixa_pmtu;               /* IXAF_VERIFY_PMTU */

2115         /* Set by ULP if IXAF_VERIFY_PMTU; otherwise set by IP */
2116         uint32_t       ixa_fragsize;

2118         int8_t         ixa_use_min_mtu;        /* IXAF_USE_MIN_MTU values */

2120         pfirepostfrag_t ixa_postfragfn;        /* Set internally in IP */

2122         in6_addr_t     ixa_nexthop_v6;         /* IXAF_NEXTHOP_SET */
2123 #define ixa_nexthop_v4  V4_PART_OF_V6(ixa_nexthop_v6)

2125         zoneid_t       ixa_no_loop_zoneid;     /* IXAF_NO_LOOP_ZONEID_SET */

2127         uint_t         ixa_scopeid;            /* For IPv6 link-locals */

2129         uint_t         ixa_broadcast_ttl;      /* IXAF_BROACAST_TTL_SET */

2131         uint_t         ixa_multicast_ttl;      /* Assumed set for multicast */
2132         uint_t         ixa_multicast_ifindex;  /* Assumed set for multicast */
2133         ipaddr_t       ixa_multicast_ifaddr;   /* Assumed set for multicast */

2135         int            ixa_raw_cksum_offset;   /* If IXAF_SET_RAW_CKSUM */

2137         uint32_t       ixa_ident;              /* For IPv6 fragment header */

2139         uint64_t       ixa_conn_id;            /* Used by DTrace */
2140         /*
2141          * Cached LSO information.
2142          */
2143         ill_lso_capab_t ixa_lso_capab;         /* Valid when IXAF_LSO_CAPAB */

2145         uint64_t       ixa_ipsec_policy_gen;   /* Generation from iph_gen */
2146         /*
2147          * The following IPsec fields are only initialized when
2148          * IXAF_IPSEC_SECURE is set. Otherwise they contain garbage.
2149          */
2150         ipsec_latch_t  *ixa_ipsec_latch;       /* Just the ids */
2151         struct ipsa_s  *ixa_ipsec_ah_sa;       /* Hard reference SA for AH */
2152         struct ipsa_s  *ixa_ipsec_esp_sa;      /* Hard reference SA for ESP */
2153         struct ipsec_policy_s  *ixa_ipsec_policy; /* why are we here? */
2154         struct ipsec_action_s  *ixa_ipsec_action; /* For reflected packets */
2155         ipsa_ref_t     ixa_ipsec_ref[2];       /* Soft reference to SA */
2156                                                /* 0: ESP, 1: AH */

2158         /*
2159          * The selectors here are potentially different than the SPD rule's
2160          * selectors, and we need to have both available for IKEv2.
2161          *
2162          * NOTE: "Source" and "Dest" are w.r.t. outbound datagrams.  Ports can
2163          *       be zero, and the protocol number is needed to make the ports
2164          *       significant.
2165          */
2166         uint16_t ixa_ipsec_src_port;   /* Source port number of d-gram. */
2167         uint16_t ixa_ipsec_dst_port;   /* Destination port number of d-gram. */
2168         uint8_t  ixa_ipsec_icmp_type;  /* ICMP type of d-gram */
2169         uint8_t  ixa_ipsec_icmp_code;  /* ICMP code of d-gram */

2171         sa_family_t ixa_ipsec_inaf;    /* Inner address family */
2172 #define IXA_MAX_ADDRLEN 4      /* Max addr len. (in 32-bit words) */
2173         uint32_t ixa_ipsec_insrc[IXA_MAX_ADDRLEN];     /* Inner src address */
```

```
2174         uint32_t ixa_ipsec_indst[IXA_MAX_ADDRLEN];     /* Inner dest address */
2175         uint8_t  ixa_ipsec_insrcpfx;   /* Inner source prefix */
2176         uint8_t  ixa_ipsec_indstpfx;   /* Inner destination prefix */

2178         uint8_t ixa_ipsec_proto;       /* IP protocol number for d-gram. */

2180         /* Always initialized independently of ixa_flags settings */
2181         uint_t         ixa_ifindex;    /* Assumed always set */
2182         uint16_t       ixa_ip_hdr_length; /* Points to ULP header */
2183         uint8_t        ixa_protocol;   /* Protocol number for ULP cksum */
2184         ts_label_t     *ixa_tsl;       /* Always set. NULL if not TX */
2185         ip_stack_t     *ixa_ipst;      /* Always set */
2186         uint32_t       ixa_extra_ident; /* Set if LSO */
2187         cred_t         *ixa_cred;      /* For getpeerucred */
2188         pid_t          ixa_cpid;       /* For getpeerucred */

2190 #ifdef DEBUG
2191         kthread_t      *ixa_curthread; /* For serialization assert */
2192 #endif
2193         squeue_t       *ixa_sqp;       /* Set from conn_sqp as a hint */
2194         uintptr_t      ixa_cookie;     /* cookie to use for tx flow control */

2196         /*
2197          * Must be set by ULP if any of IXAF_VERIFY_LSO, IXAF_VERIFY_PMTU,
2198          * or IXAF_VERIFY_ZCOPY is set.
2199          */
2200         ixa_notify_t   ixa_notify;     /* Registered upcall notify function */
2201         void           *ixa_notify_cookie; /* ULP cookie for ixa_notify */
2202 };

2204 /*
2205  * Flags to indicate which transmit attributes are set.
2206  * Split into "xxx_SET" ones which indicate that the "xxx" field it set, and
2207  * single flags.
2208  */
2209 #define IXAF_REACH_CONF         0x00000001      /* Reachability confirmation */
2210 #define IXAF_BROADCAST_TTL_SET  0x00000002      /* ixa_broadcast_ttl valid */
2211 #define IXAF_SET_SOURCE         0x00000004      /* Replace if broadcast */
2212 #define IXAF_USE_MIN_MTU        0x00000008      /* IPV6_USE_MIN_MTU */

2214 #define IXAF_DONTFRAG           0x00000010      /* IP*_DONTFRAG */
2215 #define IXAF_VERIFY_PMTU        0x00000020      /* ixa_pmtu/ixa_fragsize set */
2216 #define IXAF_PMTU_DISCOVERY     0x00000040      /* Create/use PMTU state */
2217 #define IXAF_MULTICAST_LOOP     0x00000080      /* IP_MULTICAST_LOOP */

2219 #define IXAF_IPSEC_SECURE       0x00000100      /* Need IPsec processing */
2220 #define IXAF_UCRED_TSL          0x00000200      /* ixa_tsl from SCM_UCRED */
2221 #define IXAF_DONTROUTE          0x00000400      /* SO_DONTROUTE */
2222 #define IXAF_NO_IPSEC           0x00000800      /* Ignore policy */

2224 #define IXAF_PMTU_TOO_SMALL     0x00001000      /* PMTU too small */
2225 #define IXAF_SET_ULP_CKSUM      0x00002000      /* Calculate ULP checksum */
2226 #define IXAF_VERIFY_SOURCE      0x00004000      /* Check that source is ok */
2227 #define IXAF_NEXTHOP_SET        0x00008000      /* ixa_nexthop set */

2229 #define IXAF_PMTU_IPV4_DF       0x00010000      /* Set IPv4 DF */
2230 #define IXAF_NO_DEV_FLOW_CTL    0x00020000      /* Protocol needs no flow ctl */
2231 #define IXAF_NO_TTL_CHANGE      0x00040000      /* Internal to IP */
2232 #define IXAF_IPV6_ADD_FRAGHDR   0x00080000      /* Add fragment header */

2234 #define IXAF_IPSEC_TUNNEL       0x00100000      /* Tunnel mode */
2235 #define IXAF_NO_PFHOOK          0x00200000      /* Skip xmit pfhook */
2236 #define IXAF_NO_TRACE           0x00400000      /* When back from ARP/ND */
2237 #define IXAF_SCOPEID_SET        0x00800000      /* ixa_scopeid set */

2239 #define IXAF_MULTIRT_MULTICAST  0x01000000      /* MULTIRT for multicast */
```

```
2240 #define IXAF_NO_HW_CKSUM        0x02000000      /* Force software cksum */
2241 #define IXAF_SET_RAW_CKSUM      0x04000000      /* Use ixa_raw_cksum_offset */
2242 #define IXAF_IPSEC_GLOBAL_POLICY 0x08000000     /* Policy came from global */

2244 /* Note the following uses bits 0x10000000 through 0x80000000 */
2245 #define IXAF_IS_IPV4            IAF_IS_IPV4
2246 #define IXAF_TRUSTED_ICMP       IAF_TRUSTED_ICMP
2247 #define IXAF_NO_LOOP_ZONEID_SET IAF_NO_LOOP_ZONEID_SET
2248 #define IXAF_LOOPBACK_COPY      IAF_LOOPBACK_COPY

2250 /* Note: use the upper 32 bits */
2251 #define IXAF_VERIFY_LSO         0x100000000     /* Check LSO capability */
2252 #define IXAF_LSO_CAPAB          0x200000000     /* Capable of LSO */
2253 #define IXAF_VERIFY_ZCOPY       0x400000000     /* Check Zero Copy capability */
2254 #define IXAF_ZCOPY_CAPAB        0x800000000     /* Capable of ZEROCOPY */

2256 /*
2257  * The normal flags for sending packets e.g., icmp errors
2258  */
2259 #define IXAF_BASIC_SIMPLE_V4    \
2260         (IXAF_SET_ULP_CKSUM | IXAF_IS_IPV4 | IXAF_VERIFY_SOURCE)
2261 #define IXAF_BASIC_SIMPLE_V6    (IXAF_SET_ULP_CKSUM | IXAF_VERIFY_SOURCE)

2263 /*
2264  * Normally these fields do not have a hold. But in some cases they do, for
2265  * instance when we've gone through ip_*_attr_to/from_mblk.
2266  * We use ixa_free_flags to indicate that they have a hold and need to be
2267  * released on cleanup.
2268  */
2269 #define IXA_FREE_CRED           0x00000001      /* ixa_cred needs to be rele */
2270 #define IXA_FREE_TSL            0x00000002      /* ixa_tsl needs to be rele */

2272 /*
2273  * Simplistic way to set the ixa_xmit_hint for locally generated traffic
2274  * and forwarded traffic. The shift amount are based on the size of the
2275  * structs to discard the low order bits which don't have much if any variation
2276  * (coloring in kmem_cache_alloc might provide some variation).
2277  *
2278  * Basing the locally generated hint on the address of the conn_t means that
2279  * the packets from the same socket/connection do not get reordered.
2280  * Basing the hint for forwarded traffic on the ill_ring_t means that
2281  * packets from the same NIC+ring are likely to use the same outbound ring
2282  * hence we get low contention on the ring in the transmitting driver.
2283  */
2284 #define CONN_TO_XMIT_HINT(connp)        ((uint32_t)(((uintptr_t)connp) >> 11))
2285 #define ILL_RING_TO_XMIT_HINT(ring)     ((uint32_t)(((uintptr_t)ring) >> 7))

2287 /*
2288  * IP set Destination Flags used by function ip_set_destination,
2289  * ip_attr_connect, and conn_connect.
2290  */
2291 #define IPDF_ALLOW_MCBC         0x1     /* Allow multi/broadcast */
2292 #define IPDF_VERIFY_DST         0x2     /* Verify destination addr */
2293 #define IPDF_SELECT_SRC         0x4     /* Select source address */
2294 #define IPDF_LSO                0x8     /* Try LSO */
2295 #define IPDF_IPSEC              0x10    /* Set IPsec policy */
2296 #define IPDF_ZONE_IS_GLOBAL     0x20    /* From conn_zone_is_global */
2297 #define IPDF_ZCOPY              0x40    /* Try ZEROCOPY */
2298 #define IPDF_UNIQUE_DCE         0x80    /* Get a per-destination DCE */

2300 /*
2301  * Receive side attributes used between the transport protocols and IP as
2302  * well as inside IP.
2303  */
2304 struct ip_recv_attr_s {
2305         iaflags_t       ira_flags;      /* See below */
```

```
2307         uint32_t        ira_free_flags; /* IRA_FREE_*. See below */

2309         /*
2310          * This is a hint for TCP SYN packets.
2311          * Always initialized independently of ira_flags settings
2312          */
2313         squeue_t        *ira_sqp;
2314         ill_rx_ring_t   *ira_ring;      /* Internal to IP */

2316         /* For ip_accept_tcp when IRAF_TARGET_SQP is set */
2317         squeue_t        *ira_target_sqp;
2318         mblk_t          *ira_target_sqp_mp;

2320         /* Always initialized independently of ira_flags settings */
2321         uint32_t        ira_xmit_hint; /* For ECMP and GLD TX ring fanout */
2322         zoneid_t        ira_zoneid;     /* ALL_ZONES unless local delivery */
2323         uint_t          ira_pktlen;     /* Always set. For frag and stats */
2324         uint16_t        ira_ip_hdr_length; /* Points to ULP header */
2325         uint8_t         ira_protocol;   /* Protocol number for ULP cksum */
2326         uint_t          ira_rifindex;   /* Received ifindex */
2327         uint_t          ira_ruifindex;  /* Received upper ifindex */
2328         ts_label_t      *ira_tsl;       /* Always set. NULL if not TX */
2329         /*
2330          * ira_rill and ira_ill is set inside IP, but not when conn_recv is
2331          * called; ULPs should use ira_ruifindex instead.
2332          */
2333         ill_t           *ira_rill;      /* ill where packet came */
2334         ill_t           *ira_ill;       /* ill where IP address hosted */
2335         cred_t          *ira_cred;      /* For getpeerucred */
2336         pid_t           ira_cpid;       /* For getpeerucred */

2338         /* Used when IRAF_VERIFIED_SRC is set; this source was ok */
2339         ipaddr_t        ira_verified_src;

2341         /*
2342          * The following IPsec fields are only initialized when
2343          * IRAF_IPSEC_SECURE is set. Otherwise they contain garbage.
2344          */
2345         struct ipsec_action_s *ira_ipsec_action; /* how we made it in.. */
2346         struct ipsa_s   *ira_ipsec_ah_sa;       /* SA for AH */
2347         struct ipsa_s   *ira_ipsec_esp_sa;      /* SA for ESP */

2349         ipaddr_t        ira_mroute_tunnel;      /* IRAF_MROUTE_TUNNEL_SET */

2351         zoneid_t        ira_no_loop_zoneid;     /* IRAF_NO_LOOP_ZONEID_SET */

2353         uint32_t        ira_esp_udp_ports;      /* IRAF_ESP_UDP_PORTS */

2355         /*
2356          * For IP_RECVSLLA and ip_ndp_conflict/find_solicitation.
2357          * Same size as max for sockaddr_dl
2358          */
2359 #define IRA_L2SRC_SIZE  244
2360         uint8_t         ira_l2src[IRA_L2SRC_SIZE];      /* If IRAF_L2SRC_SET */

2362         /*
2363          * Local handle that we use to do lazy setting of ira_l2src.
2364          * We defer setting l2src until needed but we do before any
2365          * ip_input pullupmsg or copymsg.
2366          */
2367         struct mac_header_info_s *ira_mhip;     /* Could be NULL */
2368 };

2370 /*
2371  * Flags to indicate which receive attributes are set.
```

```
2372  */
2373 #define IRAF_SYSTEM_LABELED      0x00000001       /* is_system_labeled() */
2374 #define IRAF_IPV4_OPTIONS        0x00000002       /* Performance */
2375 #define IRAF_MULTICAST           0x00000004       /* Was multicast at L3 */
2376 #define IRAF_BROADCAST           0x00000008       /* Was broadcast at L3 */
2377 #define IRAF_MULTIBROADCAST      (IRAF_MULTICAST|IRAF_BROADCAST)

2379 #define IRAF_LOOPBACK            0x00000010       /* Looped back by IP */
2380 #define IRAF_VERIFY_IP_CKSUM     0x00000020       /* Need to verify IP */
2381 #define IRAF_VERIFY_ULP_CKSUM    0x00000040       /* Need to verify TCP,UDP,etc */
2382 #define IRAF_SCTP_CSUM_ERR       0x00000080       /* sctp pkt has failed chksum */

2384 #define IRAF_IPSEC_SECURE        0x00000100       /* Passed AH and/or ESP */
2385 #define IRAF_DHCP_UNICAST        0x00000200
2386 #define IRAF_IPSEC_DECAPS        0x00000400       /* Was packet decapsulated */
2387                                                  /* from a matching inner packet? */
2388 #define IRAF_TARGET_SQP          0x00000800       /* ira_target_sqp is set */
2389 #define IRAF_VERIFIED_SRC        0x00001000       /* ira_verified_src set */
2390 #define IRAF_RSVP                0x00002000       /* RSVP packet for rsvpd */
2391 #define IRAF_MROUTE_TUNNEL_SET   0x00004000       /* From ip_mroute_decap */
2392 #define IRAF_PIM_REGISTER        0x00008000       /* From register_mforward */

2394 #define IRAF_TX_MAC_EXEMPTABLE   0x00010000       /* Allow MAC_EXEMPT readdown */
2395 #define IRAF_TX_SHARED_ADDR      0x00020000       /* Arrived on ALL_ZONES addr */
2396 #define IRAF_ESP_UDP_PORTS       0x00040000       /* NAT-traversal packet */
2397 #define IRAF_NO_HW_CKSUM         0x00080000       /* Force software cksum */

2399 #define IRAF_ICMP_ERROR          0x00100000       /* Send to conn_recvicmp */
2400 #define IRAF_ROUTER_ALERT        0x00200000       /* IPv6 router alert */
2401 #define IRAF_L2SRC_SET           0x00400000       /* ira_l2src has been set */
2402 #define IRAF_L2SRC_LOOPBACK      0x00800000       /* Came from us */

2404 #define IRAF_L2DST_MULTICAST     0x01000000       /* Multicast at L2 */
2405 #define IRAF_L2DST_BROADCAST     0x02000000       /* Broadcast at L2 */
2406 /* Unused 0x04000000 */
2407 /* Unused 0x08000000 */

2409 /* Below starts with 0x10000000 */
2410 #define IRAF_IS_IPV4             IAF_IS_IPV4
2411 #define IRAF_TRUSTED_ICMP        IAF_TRUSTED_ICMP
2412 #define IRAF_NO_LOOP_ZONEID_SET  IAF_NO_LOOP_ZONEID_SET
2413 #define IRAF_LOOPBACK_COPY       IAF_LOOPBACK_COPY

2415 /*
2416  * Normally these fields do not have a hold. But in some cases they do, for
2417  * instance when we've gone through ip_*_attr_to/from_mblk.
2418  * We use ira_free_flags to indicate that they have a hold and need to be
2419  * released on cleanup.
2420  */
2421 #define IRA_FREE_CRED            0x00000001       /* ira_cred needs to be rele */
2422 #define IRA_FREE_TSL             0x00000002       /* ira_tsl needs to be rele */

2424 /*
2425  * Optional destination cache entry for path MTU information,
2426  * and ULP metrics.
2427  */
2428 struct dce_s {
2429         uint_t          dce_generation; /* Changed since cached? */
2430         uint_t          dce_flags;      /* See below */
2431         uint_t          dce_ipversion;  /* IPv4/IPv6 version */
2432         uint32_t        dce_pmtu;       /* Path MTU if DCEF_PMTU */
2433         uint32_t        dce_ident;      /* Per destination IP ident. */
2434         iulp_t          dce_uinfo;      /* Metrics if DCEF_UINFO */

2436         struct dce_s    *dce_next;
2437         struct dce_s    **dce_ptpn;
```

```
2438         struct dcb_s    *dce_bucket;

2440         union {
2441                 in6_addr_t      dceu_v6addr;
2442                 ipaddr_t        dceu_v4addr;
2443         } dce_u;
2444 #define dce_v4addr       dce_u.dceu_v4addr
2445 #define dce_v6addr       dce_u.dceu_v6addr
2446         /* Note that for IPv6+IPMP we use the ifindex for the upper interface */
2447         uint_t          dce_ifindex;    /* For IPv6 link-locals */

2449         kmutex_t        dce_lock;
2450         uint_t          dce_refcnt;
2451         uint64_t        dce_last_change_time;   /* Path MTU. In seconds */

2453         ip_stack_t      *dce_ipst;       /* Does not have a netstack_hold */
2454 };

2456 /*
2457  * Values for dce_generation.
2458  *
2459  * If a DCE has DCE_GENERATION_CONDEMNED, the last dce_refrele should delete
2460  * it.
2461  *
2462  * DCE_GENERATION_VERIFY is never stored in dce_generation but it is
2463  * stored in places that cache DCE (such as ixa_dce_generation).
2464  * It is used as a signal that the cache is stale and needs to be reverified.
2465  */
2466 #define DCE_GENERATION_CONDEMNED         0
2467 #define DCE_GENERATION_VERIFY            1
2468 #define DCE_GENERATION_INITIAL           2
2469 #define DCE_IS_CONDEMNED(dce) \
2470         ((dce)->dce_generation == DCE_GENERATION_CONDEMNED)


2473 /*
2474  * Values for ips_src_generation.
2475  *
2476  * SRC_GENERATION_VERIFY is never stored in ips_src_generation but it is
2477  * stored in places that cache IREs (ixa_src_generation). It is used as a
2478  * signal that the cache is stale and needs to be reverified.
2479  */
2480 #define SRC_GENERATION_VERIFY            0
2481 #define SRC_GENERATION_INITIAL           1

2483 /*
2484  * The kernel stores security attributes of all gateways in a database made
2485  * up of one or more tsol_gcdb_t elements.  Each tsol_gcdb_t contains the
2486  * security-related credentials of the gateway.  More than one gateways may
2487  * share entries in the database.
2488  *
2489  * The tsol_gc_t structure represents the gateway to credential association,
2490  * and refers to an entry in the database.  One or more tsol_gc_t entities are
2491  * grouped together to form one or more tsol_gcgrp_t, each representing the
2492  * list of security attributes specific to the gateway.  A gateway may be
2493  * associated with at most one credentials group.
2494  */
2495 struct tsol_gcgrp_s;

2497 extern uchar_t  ip6opt_ls;      /* TX IPv6 enabler */

2499 /*
2500  * Gateway security credential record.
2501  */
2502 typedef struct tsol_gcdb_s {
2503         uint_t          gcdb_refcnt;    /* reference count */
```

```
2504          struct rtsa_s   gcdb_attr;      /* security attributes */
2505 #define gcdb_mask        gcdb_attr.rtsa_mask
2506 #define gcdb_doi         gcdb_attr.rtsa_doi
2507 #define gcdb_slrange     gcdb_attr.rtsa_slrange
2508 } tsol_gcdb_t;

2510 /*
2511  * Gateway to credential association.
2512  */
2513 typedef struct tsol_gc_s {
2514          uint_t          gc_refcnt;      /* reference count */
2515          struct tsol_gcgrp_s *gc_grp;    /* pointer to group */
2516          struct tsol_gc_s *gc_prev;      /* previous in list */
2517          struct tsol_gc_s *gc_next;      /* next in list */
2518          tsol_gcdb_t     *gc_db;         /* pointer to actual credentials */
2519 } tsol_gc_t;

2521 /*
2522  * Gateway credentials group address.
2523  */
2524 typedef struct tsol_gcgrp_addr_s {
2525          int             ga_af;          /* address family */
2526          in6_addr_t      ga_addr;        /* IPv4 mapped or IPv6 address */
2527 } tsol_gcgrp_addr_t;

2529 /*
2530  * Gateway credentials group.
2531  */
2532 typedef struct tsol_gcgrp_s {
2533          uint_t          gcgrp_refcnt;   /* reference count */
2534          krwlock_t       gcgrp_rwlock;   /* lock to protect following */
2535          uint_t          gcgrp_count;    /* number of credentials */
2536          tsol_gc_t       *gcgrp_head;    /* first credential in list */
2537          tsol_gc_t       *gcgrp_tail;    /* last credential in list */
2538          tsol_gcgrp_addr_t gcgrp_addr;   /* next-hop gateway address */
2539 } tsol_gcgrp_t;

2541 extern kmutex_t gcgrp_lock;

2543 #define GC_REFRELE(p) {                                  \
2544          ASSERT((p)->gc_grp != NULL);                    \
2545          rw_enter(&(p)->gc_grp->gcgrp_rwlock, RW_WRITER); \
2546          ASSERT((p)->gc_refcnt > 0);                     \
2547          if (--((p)->gc_refcnt) == 0)                    \
2548                  gc_inactive(p);                         \
2549          else                                            \
2550                  rw_exit(&(p)->gc_grp->gcgrp_rwlock); \
2551 }

2553 #define GCGRP_REFHOLD(p) {                               \
2554          mutex_enter(&gcgrp_lock);                       \
2555          ++((p)->gcgrp_refcnt);                          \
2556          ASSERT((p)->gcgrp_refcnt != 0);                 \
2557          mutex_exit(&gcgrp_lock);                        \
2558 }

2560 #define GCGRP_REFRELE(p) {                               \
2561          mutex_enter(&gcgrp_lock);                       \
2562          ASSERT((p)->gcgrp_refcnt > 0);                  \
2563          if (--((p)->gcgrp_refcnt) == 0)                 \
2564                  gcgrp_inactive(p);                      \
2565          ASSERT(MUTEX_HELD(&gcgrp_lock));                \
2566          mutex_exit(&gcgrp_lock);                        \
2567 }

2569 /*
```

```
2570  * IRE gateway security attributes structure, pointed to by tsol_ire_gw_secattr
2571  */
2572 struct tsol_tnrhc;

2574 struct tsol_ire_gw_secattr_s {
2575          kmutex_t        igsa_lock;      /* lock to protect following */
2576          struct tsol_tnrhc *igsa_rhc;    /* host entry for gateway */
2577          tsol_gc_t       *igsa_gc;       /* for prefix IREs */
2578 };

2580 void irb_refrele_ftable(irb_t *);

2582 extern struct kmem_cache *rt_entry_cache;

2584 typedef struct ire4 {
2585          ipaddr_t ire4_mask;             /* Mask for matching this IRE. */
2586          ipaddr_t ire4_addr;             /* Address this IRE represents. */
2587          ipaddr_t ire4_gateway_addr;     /* Gateway including for IRE_ONLINK */
2588          ipaddr_t ire4_setsrc_addr;      /* RTF_SETSRC */
2589 } ire4_t;

2591 typedef struct ire6 {
2592          in6_addr_t ire6_mask;           /* Mask for matching this IRE. */
2593          in6_addr_t ire6_addr;           /* Address this IRE represents. */
2594          in6_addr_t ire6_gateway_addr;   /* Gateway including for IRE_ONLINK */
2595          in6_addr_t ire6_setsrc_addr;    /* RTF_SETSRC */
2596 } ire6_t;

2598 typedef union ire_addr {
2599          ire6_t  ire6_u;
2600          ire4_t  ire4_u;
2601 } ire_addr_u_t;

2603 /*
2604  * Internet Routing Entry
2605  * When we have multiple identical IREs we logically add them by manipulating
2606  * ire_identical_ref and ire_delete first decrements
2607  * that and when it reaches 1 we know it is the last IRE.
2608  * "identical" is defined as being the same for:
2609  * ire_addr, ire_netmask, ire_gateway, ire_ill, ire_zoneid, and ire_type
2610  * For instance, multiple IRE_BROADCASTs for the same subnet number are
2611  * viewed as identical, and so are the IRE_INTERFACEs when there are
2612  * multiple logical interfaces (on the same ill) with the same subnet prefix.
2613  */
2614 struct ire_s {
2615          struct  ire_s   *ire_next;      /* The hash chain must be first. */
2616          struct  ire_s   **ire_ptpn;     /* Pointer to previous next. */
2617          uint32_t        ire_refcnt;     /* Number of references */
2618          ill_t           *ire_ill;
2619          uint32_t        ire_identical_ref; /* IRE_INTERFACE, IRE_BROADCAST */
2620          uchar_t         ire_ipversion;  /* IPv4/IPv6 version */
2621          ushort_t        ire_type;       /* Type of IRE */
2622          uint_t          ire_generation; /* Generation including CONDEMNED */
2623          uint_t  ire_ib_pkt_count;       /* Inbound packets for ire_addr */
2624          uint_t  ire_ob_pkt_count;       /* Outbound packets to ire_addr */
2625          time_t  ire_create_time;        /* Time (in secs) IRE was created. */
2626          uint32_t        ire_flags;      /* flags related to route (RTF_*) */
2627          /*
2628           * ire_testhidden is TRUE for INTERFACE IREs of IS_UNDER_IPMP(ill)
2629           * interfaces
2630           */
2631          boolean_t       ire_testhidden;
2632          pfirerecv_t     ire_recvfn;     /* Receive side handling */
2633          pfiresend_t     ire_sendfn;     /* Send side handling */
2634          pfirepostfrag_t ire_postfragfn; /* Bottom end of send handling */
```

```
2636            uint_t          ire_masklen;    /* # bits in ire_mask{,_v6} */
2637            ire_addr_u_t    ire_u;          /* IPv4/IPv6 address info. */

2639            irb_t           *ire_bucket;    /* Hash bucket when ire_ptphn is set */
2640            kmutex_t        ire_lock;
2641            clock_t         ire_last_used_time;     /* For IRE_LOCAL reception */
2642            tsol_ire_gw_secattr_t *ire_gw_secattr; /* gateway security attributes */
2643            zoneid_t        ire_zoneid;

2645            /*
2646             * Cached information of where to send packets that match this route.
2647             * The ire_dep_* information is used to determine when ire_nce_cache
2648             * needs to be updated.
2649             * ire_nce_cache is the fastpath for the Neighbor Cache Entry
2650             * for IPv6; arp info for IPv4
2651             * Since this is a cache setup and torn down independently of
2652             * applications we need to use nce_ref{rele,hold}_notr for it.
2653             */
2654            nce_t           *ire_nce_cache;

2656            /*
2657             * Quick check whether the ire_type and ire_masklen indicates
2658             * that the IRE can have ire_nce_cache set i.e., whether it is
2659             * IRE_ONLINK and for a single destination.
2660             */
2661            boolean_t       ire_nce_capable;

2663            /*
2664             * Dependency tracking so we can safely cache IRE and NCE pointers
2665             * in offlink and onlink IREs.
2666             * These are locked under the ips_ire_dep_lock rwlock. Write held
2667             * when modifying the linkage.
2668             * ire_dep_parent (Also chain towards IRE for nexthop)
2669             * ire_dep_parent_generation: ire_generation of ire_dep_parent
2670             * ire_dep_children (From parent to first child)
2671             * ire_dep_sib_next (linked list of siblings)
2672             * ire_dep_sib_ptpn (linked list of siblings)
2673             *
2674             * The parent has a ire_refhold on each child, and each child has
2675             * an ire_refhold on its parent.
2676             * Since ire_dep_parent is a cache setup and torn down independently of
2677             * applications we need to use ire_ref{rele,hold}_notr for it.
2678             */
2679            ire_t           *ire_dep_parent;
2680            ire_t           *ire_dep_children;
2681            ire_t           *ire_dep_sib_next;
2682            ire_t           **ire_dep_sib_ptpn;    /* Pointer to previous next */
2683            uint_t          ire_dep_parent_generation;

2685            uint_t          ire_badcnt;     /* Number of times ND_UNREACHABLE */
2686            uint64_t        ire_last_badcnt;       /* In seconds */

2688            /* ire_defense* and ire_last_used_time are only used on IRE_LOCALs */
2689            uint_t          ire_defense_count;     /* number of ARP conflicts */
2690            uint_t          ire_defense_time;      /* last time defended (secs) */

2692            boolean_t       ire_trace_disable;     /* True when alloc fails */
2693            ip_stack_t      *ire_ipst;      /* Does not have a netstack_hold */
2694            iulp_t          ire_metrics;
2695            /*
2696             * default and prefix routes that are added without explicitly
2697             * specifying the interface are termed "unbound" routes, and will
2698             * have ire_unbound set to true.
2699             */
2700            boolean_t       ire_unbound;
2701 };
```

```
2703 /* IPv4 compatibility macros */
2704 #define ire_mask                ire_u.ire4_u.ire4_mask
2705 #define ire_addr                ire_u.ire4_u.ire4_addr
2706 #define ire_gateway_addr        ire_u.ire4_u.ire4_gateway_addr
2707 #define ire_setsrc_addr         ire_u.ire4_u.ire4_setsrc_addr

2709 #define ire_mask_v6             ire_u.ire6_u.ire6_mask
2710 #define ire_addr_v6             ire_u.ire6_u.ire6_addr
2711 #define ire_gateway_addr_v6     ire_u.ire6_u.ire6_gateway_addr
2712 #define ire_setsrc_addr_v6      ire_u.ire6_u.ire6_setsrc_addr

2714 /*
2715  * Values for ire_generation.
2716  *
2717  * If an IRE is marked with IRE_IS_CONDEMNED, the last walker of
2718  * the bucket should delete this IRE from this bucket.
2719  *
2720  * IRE_GENERATION_VERIFY is never stored in ire_generation but it is
2721  * stored in places that cache IREs (such as ixa_ire_generation and
2722  * ire_dep_parent_generation). It is used as a signal that the cache is
2723  * stale and needs to be reverified.
2724  */
2725 #define IRE_GENERATION_CONDEMNED        0
2726 #define IRE_GENERATION_VERIFY           1
2727 #define IRE_GENERATION_INITIAL          2
2728 #define IRE_IS_CONDEMNED(ire) \
2729         ((ire)->ire_generation == IRE_GENERATION_CONDEMNED)

2731 /* Convenient typedefs for sockaddrs */
2732 typedef struct sockaddr_in      sin_t;
2733 typedef struct sockaddr_in6     sin6_t;

2735 /* Name/Value Descriptor. */
2736 typedef struct nv_s {
2737         uint64_t nv_value;
2738         char    *nv_name;
2739 } nv_t;

2741 #define ILL_FRAG_HASH(s, i) \
2742         ((ntohl(s) ^ ((i) ^ ((i) >> 8))) % ILL_FRAG_HASH_TBL_COUNT)

2744 /*
2745  * The MAX number of allowed fragmented packets per hash bucket
2746  * calculation is based on the most common mtu size of 1500. This limit
2747  * will work well for other mtu sizes as well.
2748  */
2749 #define COMMON_IP_MTU 1500
2750 #define MAX_FRAG_MIN 10
2751 #define MAX_FRAG_PKTS(ipst)     \
2752         MAX(MAX_FRAG_MIN, (2 * (ipst->ips_ip_reass_queue_bytes / \
2753                 (COMMON_IP_MTU * ILL_FRAG_HASH_TBL_COUNT))))

2755 /*
2756  * Maximum dups allowed per packet.
2757  */
2758 extern uint_t ip_max_frag_dups;

2760 /*
2761  * Per-packet information for received packets and transmitted.
2762  * Used by the transport protocols when converting between the packet
2763  * and ancillary data and socket options.
2764  *
2765  * Note: This private data structure and related IPPF_* constant
2766  * definitions are exposed to enable compilation of some debugging tools
2767  * like lsof which use struct tcp_t in <inet/tcp.h>. This is intended to be
```

```
2768     * a temporary hack and long term alternate interfaces should be defined
2769     * to support the needs of such tools and private definitions moved to
2770     * private headers.
2771     */
2772    struct ip_pkt_s {
2773            uint_t          ipp_fields;             /* Which fields are valid */
2774            in6_addr_t      ipp_addr;               /* pktinfo src/dst addr */
2775    #define ipp_addr_v4     V4_PART_OF_V6(ipp_addr)
2776            uint_t          ipp_unicast_hops;       /* IPV6_UNICAST_HOPS, IP_TTL */
2777            uint_t          ipp_hoplimit;           /* IPV6_HOPLIMIT */
2778            uint_t          ipp_hopoptslen;
2779            uint_t          ipp_rthdrdstoptslen;
2780            uint_t          ipp_rthdrlen;
2781            uint_t          ipp_dstoptslen;
2782            uint_t          ipp_fraghdrlen;
2783            ip6_hbh_t       *ipp_hopopts;
2784            ip6_dest_t      *ipp_rthdrdstopts;
2785            ip6_rthdr_t     *ipp_rthdr;
2786            ip6_dest_t      *ipp_dstopts;
2787            ip6_frag_t      *ipp_fraghdr;
2788            uint8_t         ipp_tclass;             /* IPV6_TCLASS */
2789            uint8_t         ipp_type_of_service;    /* IP_TOS */
2790            uint_t          ipp_ipv4_options_len;   /* Len of IPv4 options */
2791            uint8_t         *ipp_ipv4_options;      /* Ptr to IPv4 options */
2792            uint_t          ipp_label_len_v4;       /* Len of TX label for IPv4 */
2793            uint8_t         *ipp_label_v4;          /* TX label for IPv4 */
2794            uint_t          ipp_label_len_v6;       /* Len of TX label for IPv6 */
2795            uint8_t         *ipp_label_v6;          /* TX label for IPv6 */
2796    };
2797    typedef struct ip_pkt_s ip_pkt_t;

2799    extern void ip_pkt_free(ip_pkt_t *);    /* free storage inside ip_pkt_t */
2800    extern ipaddr_t ip_pkt_source_route_v4(const ip_pkt_t *);
2801    extern in6_addr_t *ip_pkt_source_route_v6(const ip_pkt_t *);
2802    extern int ip_pkt_copy(ip_pkt_t *, ip_pkt_t *, int);
2803    extern void ip_pkt_source_route_reverse_v4(ip_pkt_t *);

2805    /* ipp_fields values */
2806    #define IPPF_ADDR               0x0001  /* Part of in6_pktinfo: src/dst addr */
2807    #define IPPF_HOPLIMIT           0x0002  /* Overrides unicast and multicast */
2808    #define IPPF_TCLASS             0x0004  /* Overrides class in sin6_flowinfo */

2810    #define IPPF_HOPOPTS            0x0010  /* ipp_hopopts set */
2811    #define IPPF_RTHDR              0x0020  /* ipp_rthdr set */
2812    #define IPPF_RTHDRDSTOPTS       0x0040  /* ipp_rthdrdstopts set */
2813    #define IPPF_DSTOPTS            0x0080  /* ipp_dstopts set */

2815    #define IPPF_IPV4_OPTIONS       0x0100  /* ipp_ipv4_options set */
2816    #define IPPF_LABEL_V4           0x0200  /* ipp_label_v4 set */
2817    #define IPPF_LABEL_V6           0x0400  /* ipp_label_v6 set */

2819    #define IPPF_FRAGHDR            0x0800  /* Used for IPsec receive side */

2821    /*
2822     * Data structure which is passed to conn_opt_get/set.
2823     * The conn_t is included even though it can be inferred from queue_t.
2824     * setsockopt and getsockopt use conn_ixa and conn_xmit_ipp. However,
2825     * when handling ancillary data we use separate ixa and ipps.
2826     */
2827    typedef struct conn_opt_arg_s {
2828            conn_t          *coa_connp;
2829            ip_xmit_attr_t  *coa_ixa;
2830            ip_pkt_t        *coa_ipp;
2831            boolean_t       coa_ancillary;  /* Ancillary data and not setsockopt */
2832            uint_t          coa_changed;    /* See below */
2833    } conn_opt_arg_t;
```

```
2835    /*
2836     * Flags for what changed.
2837     * If we want to be more efficient in the future we can have more fine
2838     * grained flags e.g., a flag for just IP_TOS changing.
2839     * For now we either call ip_set_destination (for "route changed")
2840     * and/or conn_build_hdr_template/conn_prepend_hdr (for "header changed").
2841     */
2842    #define COA_HEADER_CHANGED      0x0001
2843    #define COA_ROUTE_CHANGED       0x0002
2844    #define COA_RCVBUF_CHANGED      0x0004  /* SO_RCVBUF */
2845    #define COA_SNDBUF_CHANGED      0x0008  /* SO_SNDBUF */
2846    #define COA_WROFF_CHANGED       0x0010  /* Header size changed */
2847    #define COA_ICMP_BIND_NEEDED    0x0020
2848    #define COA_OOBINLINE_CHANGED   0x0040

2850    #define TCP_PORTS_OFFSET        0
2851    #define UDP_PORTS_OFFSET        0

2853    /*
2854     * lookups return the ill/ipif only if the flags are clear OR Iam writer.
2855     * ill / ipif lookup functions increment the refcnt on the ill / ipif only
2856     * after calling these macros. This ensures that the refcnt on the ipif or
2857     * ill will eventually drop down to zero.
2858     */
2859    #define ILL_LOOKUP_FAILED       1       /* Used as error code */
2860    #define IPIF_LOOKUP_FAILED      2       /* Used as error code */

2862    #define ILL_CAN_LOOKUP(ill)                                             \
2863            (!((ill)->ill_state_flags & ILL_CONDEMNED) ||                   \
2864            IAM_WRITER_ILL(ill))

2866    #define ILL_IS_CONDEMNED(ill)   \
2867            ((ill)->ill_state_flags & ILL_CONDEMNED)

2869    #define IPIF_CAN_LOOKUP(ipif)   \
2870            (!((ipif)->ipif_state_flags & IPIF_CONDEMNED) || \
2871            IAM_WRITER_IPIF(ipif))

2873    #define IPIF_IS_CONDEMNED(ipif) \
2874            ((ipif)->ipif_state_flags & IPIF_CONDEMNED)

2876    #define IPIF_IS_CHANGING(ipif)  \
2877            ((ipif)->ipif_state_flags & IPIF_CHANGING)

2879    /* Macros used to assert that this thread is a writer */
2880    #define IAM_WRITER_IPSQ(ipsq)   ((ipsq)->ipsq_xop->ipx_writer == curthread)
2881    #define IAM_WRITER_ILL(ill)     IAM_WRITER_IPSQ((ill)->ill_phyint->phyint_ipsq)
2882    #define IAM_WRITER_IPIF(ipif)   IAM_WRITER_ILL((ipif)->ipif_ill)

2884    /*
2885     * Grab ill locks in the proper order. The order is highest addressed
2886     * ill is locked first.
2887     */
2888    #define GRAB_ILL_LOCKS(ill_1, ill_2)                                    \
2889    {                                                                       \
2890            if ((ill_1) > (ill_2)) {                                        \
2891                    if (ill_1 != NULL)                                      \
2892                            mutex_enter(&(ill_1)->ill_lock);                \
2893                    if (ill_2 != NULL)                                      \
2894                            mutex_enter(&(ill_2)->ill_lock);                \
2895            } else {                                                        \
2896                    if (ill_2 != NULL)                                      \
2897                            mutex_enter(&(ill_2)->ill_lock);                \
2898                    if (ill_1 != NULL && ill_1 != ill_2)                    \
2899                            mutex_enter(&(ill_1)->ill_lock);                \
```

```
2900             }                                                                 \
2901 }

2903 #define RELEASE_ILL_LOCKS(ill_1, ill_2)             \
2904 {                                                    \
2905             if (ill_1 != NULL)                       \
2906                     mutex_exit(&(ill_1)->ill_lock); \
2907             if (ill_2 != NULL && ill_2 != ill_1)     \
2908                     mutex_exit(&(ill_2)->ill_lock); \
2909 }

2911 /* Get the other protocol instance ill */
2912 #define ILL_OTHER(ill)                                               \
2913             ((ill)->ill_isv6 ? (ill)->ill_phyint->phyint_illv4 :     \
2914                 (ill)->ill_phyint->phyint_illv6)

2916 /* ioctl command info: Ioctl properties extracted and stored in here */
2917 typedef struct cmd_info_s
2918 {
2919             ipif_t  *ci_ipif;       /* ipif associated with [l]ifreq ioctl's */
2920             sin_t   *ci_sin;        /* the sin struct passed down */
2921             sin6_t  *ci_sin6;       /* the sin6_t struct passed down */
2922             struct lifreq *ci_lifr; /* the lifreq struct passed down */
2923 } cmd_info_t;

2925 extern struct kmem_cache *ire_cache;

2927 extern ipaddr_t ip_g_all_ones;

2929 extern uint_t   ip_loopback_mtu;         /* /etc/system */
2930 extern uint_t   ip_loopback_mtuplus;
2931 extern uint_t   ip_loopback_mtu_v6plus;

2933 extern vmem_t *ip_minor_arena_sa;
2934 extern vmem_t *ip_minor_arena_la;

2936 /*
2937  * ip_g_forward controls IP forwarding.  It takes two values:
2938  *      0: IP_FORWARD_NEVER    Don't forward packets ever.
2939  *      1: IP_FORWARD_ALWAYS   Forward packets for elsewhere.
2940  *
2941  * RFC1122 says there must be a configuration switch to control forwarding,
2942  * but the default MUST be to not forward packets ever.  Implicit
2943  * control based on configuration of multiple interfaces MUST NOT be
2944  * implemented (Section 3.1).  SunOS 4.1 did provide the "automatic" capability
2945  * and, in fact, it was the default.  That capability is now provided in the
2946  * /etc/rc2.d/S69inet script.
2947  */

2949 #define ips_ip_respond_to_address_mask_broadcast \
2950                                       ips_propinfo_tbl[0].prop_cur_bval
2951 #define ips_ip_g_resp_to_echo_bcast       ips_propinfo_tbl[1].prop_cur_bval
2952 #define ips_ip_g_resp_to_echo_mcast       ips_propinfo_tbl[2].prop_cur_bval
2953 #define ips_ip_g_resp_to_timestamp        ips_propinfo_tbl[3].prop_cur_bval
2954 #define ips_ip_g_resp_to_timestamp_bcast ips_propinfo_tbl[4].prop_cur_bval
2955 #define ips_ip_g_send_redirects           ips_propinfo_tbl[5].prop_cur_bval
2956 #define ips_ip_g_forward_directed_bcast ips_propinfo_tbl[6].prop_cur_bval
2957 #define ips_ip_mrtdebug                   ips_propinfo_tbl[7].prop_cur_uval
2958 #define ips_ip_ire_reclaim_fraction       ips_propinfo_tbl[8].prop_cur_uval
2959 #define ips_ip_nce_reclaim_fraction       ips_propinfo_tbl[9].prop_cur_uval
2960 #define ips_ip_dce_reclaim_fraction       ips_propinfo_tbl[10].prop_cur_uval
2961 #define ips_ip_def_ttl                    ips_propinfo_tbl[11].prop_cur_uval
2962 #define ips_ip_forward_src_routed         ips_propinfo_tbl[12].prop_cur_bval
2963 #define ips_ip_wroff_extra                ips_propinfo_tbl[13].prop_cur_uval
2964 #define ips_ip_pathmtu_interval           ips_propinfo_tbl[14].prop_cur_uval
2965 #define ips_ip_icmp_return                ips_propinfo_tbl[15].prop_cur_uval
```

```
2966 #define ips_ip_path_mtu_discovery         ips_propinfo_tbl[16].prop_cur_bval
2967 #define ips_ip_pmtu_min                   ips_propinfo_tbl[17].prop_cur_uval
2968 #define ips_ip_ignore_redirect            ips_propinfo_tbl[18].prop_cur_bval
2969 #define ips_ip_arp_icmp_error             ips_propinfo_tbl[19].prop_cur_bval
2970 #define ips_ip_broadcast_ttl              ips_propinfo_tbl[20].prop_cur_uval
2971 #define ips_ip_icmp_err_interval          ips_propinfo_tbl[21].prop_cur_uval
2972 #define ips_ip_icmp_err_burst             ips_propinfo_tbl[22].prop_cur_uval
2973 #define ips_ip_reass_queue_bytes          ips_propinfo_tbl[23].prop_cur_uval
2974 #define ips_ip_strict_dst_multihoming     ips_propinfo_tbl[24].prop_cur_uval
2975 #define ips_ip_addrs_per_if               ips_propinfo_tbl[25].prop_cur_uval
2976 #define ips_ipsec_override_persocket_policy ips_propinfo_tbl[26].prop_cur_bval
2977 #define ips_icmp_accept_clear_messages    ips_propinfo_tbl[27].prop_cur_bval
2978 #define ips_igmp_accept_clear_messages    ips_propinfo_tbl[28].prop_cur_bval

2980 /* IPv6 configuration knobs */
2981 #define ips_ip_delay_first_probe_time     ips_propinfo_tbl[29].prop_cur_uval
2982 #define ips_max_unicast_solicit           ips_propinfo_tbl[30].prop_cur_uval
2983 #define ips_ipv6_def_hops                 ips_propinfo_tbl[31].prop_cur_uval
2984 #define ips_ipv6_icmp_return              ips_propinfo_tbl[32].prop_cur_uval
2985 #define ips_ipv6_forward_src_routed       ips_propinfo_tbl[33].prop_cur_bval
2986 #define ips_ipv6_resp_echo_mcast          ips_propinfo_tbl[34].prop_cur_bval
2987 #define ips_ipv6_send_redirects           ips_propinfo_tbl[35].prop_cur_bval
2988 #define ips_ipv6_ignore_redirect          ips_propinfo_tbl[36].prop_cur_bval
2989 #define ips_ipv6_strict_dst_multihoming ips_propinfo_tbl[37].prop_cur_uval
2990 #define ips_src_check                     ips_propinfo_tbl[38].prop_cur_uval
2991 #define ips_ipsec_policy_log_interval     ips_propinfo_tbl[39].prop_cur_uval
2992 #define ips_pim_accept_clear_messages     ips_propinfo_tbl[40].prop_cur_bval
2993 #define ips_ip_ndp_unsolicit_interval     ips_propinfo_tbl[41].prop_cur_uval
2994 #define ips_ip_ndp_unsolicit_count        ips_propinfo_tbl[42].prop_cur_uval
2995 #define ips_ipv6_ignore_home_address_opt ips_propinfo_tbl[43].prop_cur_bval

2997 /* Misc IP configuration knobs */
2998 #define ips_ip_policy_mask                ips_propinfo_tbl[44].prop_cur_uval
2999 #define ips_ip_ecmp_behavior              ips_propinfo_tbl[45].prop_cur_uval
3000 #define ips_ip_multirt_ttl                ips_propinfo_tbl[46].prop_cur_uval
3001 #define ips_ip_ire_badcnt_lifetime        ips_propinfo_tbl[47].prop_cur_uval
3002 #define ips_ip_max_temp_idle              ips_propinfo_tbl[48].prop_cur_uval
3003 #define ips_ip_max_temp_defend            ips_propinfo_tbl[49].prop_cur_uval
3004 #define ips_ip_max_defend                 ips_propinfo_tbl[50].prop_cur_uval
3005 #define ips_ip_defend_interval            ips_propinfo_tbl[51].prop_cur_uval
3006 #define ips_ip_dup_recovery               ips_propinfo_tbl[52].prop_cur_uval
3007 #define ips_ip_restrict_interzone_loopback ips_propinfo_tbl[53].prop_cur_bval
3008 #define ips_ip_lso_outbound               ips_propinfo_tbl[54].prop_cur_bval
3009 #define ips_igmp_max_version              ips_propinfo_tbl[55].prop_cur_uval
3010 #define ips_mld_max_version               ips_propinfo_tbl[56].prop_cur_uval
3011 #define ips_ip_forwarding                 ips_propinfo_tbl[57].prop_cur_bval
3012 #define ips_ipv6_forwarding               ips_propinfo_tbl[58].prop_cur_bval
3013 #define ips_ip_reassembly_timeout         ips_propinfo_tbl[59].prop_cur_uval
3014 #define ips_ipv6_reassembly_timeout       ips_propinfo_tbl[60].prop_cur_uval
3015 #define ips_ip_cgtp_filter                ips_propinfo_tbl[61].prop_cur_bval
3016 #define ips_arp_probe_delay               ips_propinfo_tbl[62].prop_cur_uval
3017 #define ips_arp_fastprobe_delay           ips_propinfo_tbl[63].prop_cur_uval
3018 #define ips_arp_probe_interval            ips_propinfo_tbl[64].prop_cur_uval
3019 #define ips_arp_fastprobe_interval        ips_propinfo_tbl[65].prop_cur_uval
3020 #define ips_arp_probe_count               ips_propinfo_tbl[66].prop_cur_uval
3021 #define ips_arp_fastprobe_count           ips_propinfo_tbl[67].prop_cur_uval
3022 #define ips_ipv4_dad_announce_interval    ips_propinfo_tbl[68].prop_cur_uval
3023 #define ips_ipv6_dad_announce_interval    ips_propinfo_tbl[69].prop_cur_uval
3024 #define ips_arp_defend_interval           ips_propinfo_tbl[70].prop_cur_uval
3025 #define ips_arp_defend_rate               ips_propinfo_tbl[71].prop_cur_uval
3026 #define ips_ndp_defend_interval           ips_propinfo_tbl[72].prop_cur_uval
3027 #define ips_ndp_defend_rate               ips_propinfo_tbl[73].prop_cur_uval
3028 #define ips_arp_defend_period             ips_propinfo_tbl[74].prop_cur_uval
3029 #define ips_ndp_defend_period             ips_propinfo_tbl[75].prop_cur_uval
3030 #define ips_ipv4_icmp_return_pmtu         ips_propinfo_tbl[76].prop_cur_bval
3031 #define ips_ipv6_icmp_return_pmtu         ips_propinfo_tbl[77].prop_cur_bval
```

```
3032 #define ips_ip_arp_publish_count        ips_propinfo_tbl[78].prop_cur_uval
3033 #define ips_ip_arp_publish_interval     ips_propinfo_tbl[79].prop_cur_uval
3034 #define ips_ip_strict_src_multihoming   ips_propinfo_tbl[80].prop_cur_uval
3035 #define ips_ipv6_strict_src_multihoming ips_propinfo_tbl[81].prop_cur_uval
3036 #define ips_ipv6_drop_inbound_icmpv6    ips_propinfo_tbl[82].prop_cur_bval

3038 extern int      dohwcksum;      /* use h/w cksum if supported by the h/w */
3039 #ifdef ZC_TEST
3040 extern int      noswcksum;
3041 #endif

3043 extern char     ipif_loopback_name[];

3045 extern nv_t     *ire_nv_tbl;

3047 extern struct module_info ip_mod_info;

3049 #define HOOKS4_INTERESTED_PHYSICAL_IN(ipst)     \
3050         ((ipst)->ips_ip4_physical_in_event.he_interested)
3051 #define HOOKS6_INTERESTED_PHYSICAL_IN(ipst)     \
3052         ((ipst)->ips_ip6_physical_in_event.he_interested)
3053 #define HOOKS4_INTERESTED_PHYSICAL_OUT(ipst)    \
3054         ((ipst)->ips_ip4_physical_out_event.he_interested)
3055 #define HOOKS6_INTERESTED_PHYSICAL_OUT(ipst)    \
3056         ((ipst)->ips_ip6_physical_out_event.he_interested)
3057 #define HOOKS4_INTERESTED_FORWARDING(ipst)      \
3058         ((ipst)->ips_ip4_forwarding_event.he_interested)
3059 #define HOOKS6_INTERESTED_FORWARDING(ipst)      \
3060         ((ipst)->ips_ip6_forwarding_event.he_interested)
3061 #define HOOKS4_INTERESTED_LOOPBACK_IN(ipst)     \
3062         ((ipst)->ips_ip4_loopback_in_event.he_interested)
3063 #define HOOKS6_INTERESTED_LOOPBACK_IN(ipst)     \
3064         ((ipst)->ips_ip6_loopback_in_event.he_interested)
3065 #define HOOKS4_INTERESTED_LOOPBACK_OUT(ipst)    \
3066         ((ipst)->ips_ip4_loopback_out_event.he_interested)
3067 #define HOOKS6_INTERESTED_LOOPBACK_OUT(ipst)    \
3068         ((ipst)->ips_ip6_loopback_out_event.he_interested)
3069 /*
3070  * Hooks marcos used inside of ip
3071  * The callers use the above INTERESTED macros first, hence
3072  * the he_interested check is superflous.
3073  */
3074 #define FW_HOOKS(_hook, _event, _ilp, _olp, _iph, _fm, _m, _llm, ipst, _err) \
3075         if ((_hook).he_interested) {                                    \
3076                 hook_pkt_event_t info;                                  \
3077                                                                         \
3078                 _NOTE(CONSTCOND)                                        \
3079                 ASSERT((_ilp != NULL) || (_olp != NULL));               \
3080                                                                         \
3081                 FW_SET_ILL_INDEX(info.hpe_ifp, (ill_t *)_ilp);          \
3082                 FW_SET_ILL_INDEX(info.hpe_ofp, (ill_t *)_olp);          \
3083                 info.hpe_protocol = ipst->ips_ipv4_net_data;            \
3084                 info.hpe_hdr = _iph;                                    \
3085                 info.hpe_mp = &(_fm);                                   \
3086                 info.hpe_mb = _m;                                       \
3087                 info.hpe_flags = _llm;                                  \
3088                 _err = hook_run(ipst->ips_ipv4_net_data->netd_hooks,    \
3089                     _event, (hook_data_t)&info);                        \
3090                 if (_err != 0) {                                        \
3091                         ip2dbg(("%s hook dropped mblk chain %p hdr %p\n",\
3092                             (_hook).he_name, (void *)_fm, (void *)_m)); \
3093                         if (_fm != NULL) {                              \
3094                                 freemsg(_fm);                           \
3095                                 _fm = NULL;                             \
3096                         }                                               \
3097                         _iph = NULL;                                    \
```

```
3098                         _m = NULL;                                      \
3099                 } else {                                                \
3100                         _iph = info.hpe_hdr;                            \
3101                         _m = info.hpe_mb;                               \
3102                 }                                                       \
3103         }

3105 #define FW_HOOKS6(_hook, _event, _ilp, _olp, _iph, _fm, _m, _llm, ipst, _err) \
3106         if ((_hook).he_interested) {                                    \
3107                 hook_pkt_event_t info;                                  \
3108                                                                         \
3109                 _NOTE(CONSTCOND)                                        \
3110                 ASSERT((_ilp != NULL) || (_olp != NULL));               \
3111                                                                         \
3112                 FW_SET_ILL_INDEX(info.hpe_ifp, (ill_t *)_ilp);          \
3113                 FW_SET_ILL_INDEX(info.hpe_ofp, (ill_t *)_olp);          \
3114                 info.hpe_protocol = ipst->ips_ipv6_net_data;            \
3115                 info.hpe_hdr = _iph;                                    \
3116                 info.hpe_mp = &(_fm);                                   \
3117                 info.hpe_mb = _m;                                       \
3118                 info.hpe_flags = _llm;                                  \
3119                 _err = hook_run(ipst->ips_ipv6_net_data->netd_hooks,    \
3120                     _event, (hook_data_t)&info);                        \
3121                 if (_err != 0) {                                        \
3122                         ip2dbg(("%s hook dropped mblk chain %p hdr %p\n",\
3123                             (_hook).he_name, (void *)_fm, (void *)_m)); \
3124                         if (_fm != NULL) {                              \
3125                                 freemsg(_fm);                           \
3126                                 _fm = NULL;                             \
3127                         }                                               \
3128                         _iph = NULL;                                    \
3129                         _m = NULL;                                      \
3130                 } else {                                                \
3131                         _iph = info.hpe_hdr;                            \
3132                         _m = info.hpe_mb;                               \
3133                 }                                                       \
3134         }

3136 #define FW_SET_ILL_INDEX(fp, ill)                                       \
3137         _NOTE(CONSTCOND)                                                \
3138         if ((ill) == NULL || (ill)->ill_phyint == NULL) {               \
3139                 (fp) = 0;                                               \
3140                 _NOTE(CONSTCOND)                                        \
3141         } else if (IS_UNDER_IPMP(ill)) {                                \
3142                 (fp) = ipmp_ill_get_ipmp_ifindex(ill);                  \
3143         } else {                                                        \
3144                 (fp) = (ill)->ill_phyint->phyint_ifindex;               \
3145         }

3147 /*
3148  * Network byte order macros
3149  */
3150 #ifdef _BIG_ENDIAN
3151 #define N_IN_CLASSA_NET         IN_CLASSA_NET
3152 #define N_IN_CLASSD_NET         IN_CLASSD_NET
3153 #define N_INADDR_UNSPEC_GROUP   INADDR_UNSPEC_GROUP
3154 #define N_IN_LOOPBACK_NET       (ipaddr_t)0x7f000000U
3155 #else /* _BIG_ENDIAN */
3156 #define N_IN_CLASSA_NET         (ipaddr_t)0x000000ffU
3157 #define N_IN_CLASSD_NET         (ipaddr_t)0x000000f0U
3158 #define N_INADDR_UNSPEC_GROUP   (ipaddr_t)0x000000e0U
3159 #define N_IN_LOOPBACK_NET       (ipaddr_t)0x0000007fU
3160 #endif /* _BIG_ENDIAN */
3161 #define CLASSD(addr)    (((addr) & N_IN_CLASSD_NET) == N_INADDR_UNSPEC_GROUP)
3162 #define CLASSE(addr)    (((addr) & N_IN_CLASSD_NET) == N_IN_CLASSD_NET)
3163 #define IP_LOOPBACK_ADDR(addr)                                          \
```

```
3164         (((addr) & N_IN_CLASSA_NET == N_IN_LOOPBACK_NET))

3166 extern int       ip_debug;
3167 extern uint_t    ip_thread_data;
3168 extern krwlock_t ip_thread_rwlock;
3169 extern list_t    ip_thread_list;

3171 #ifdef IP_DEBUG
3172 #include <sys/debug.h>
3173 #include <sys/promif.h>

3175 #define ip0dbg(a)        printf a
3176 #define ip1dbg(a)        if (ip_debug > 2) printf a
3177 #define ip2dbg(a)        if (ip_debug > 3) printf a
3178 #define ip3dbg(a)        if (ip_debug > 4) printf a
3179 #else
3180 #define ip0dbg(a)        /* */
3181 #define ip1dbg(a)        /* */
3182 #define ip2dbg(a)        /* */
3183 #define ip3dbg(a)        /* */
3184 #endif  /* IP_DEBUG */

3186 /* Default MAC-layer address string length for mac_colon_addr */
3187 #define MAC_STR_LEN     128

3189 struct  mac_header_info_s;

3191 extern void     ill_frag_timer(void *);
3192 extern ill_t    *ill_first(int, int, ill_walk_context_t *, ip_stack_t *);
3193 extern ill_t    *ill_next(ill_walk_context_t *, ill_t *);
3194 extern void     ill_frag_timer_start(ill_t *);
3195 extern void     ill_nic_event_dispatch(ill_t *, lif_if_t, nic_event_t,
3196     nic_event_data_t, size_t);
3197 extern mblk_t   *ip_carve_mp(mblk_t **, ssize_t);
3198 extern mblk_t   *ip_dlpi_alloc(size_t, t_uscalar_t);
3199 extern mblk_t   *ip_dlnotify_alloc(uint_t, uint_t);
3200 extern mblk_t   *ip_dlnotify_alloc2(uint_t, uint_t, uint_t);
3201 extern char     *ip_dot_addr(ipaddr_t, char *);
3202 extern const char *mac_colon_addr(const uint8_t *, size_t, char *, size_t);
3203 extern void     ip_lwput(queue_t *, mblk_t *);
3204 extern boolean_t icmp_err_rate_limit(ip_stack_t *);
3205 extern void     icmp_frag_needed(mblk_t *, int, ip_recv_attr_t *);
3206 extern mblk_t   *icmp_inbound_v4(mblk_t *, ip_recv_attr_t *);
3207 extern void     icmp_time_exceeded(mblk_t *, uint8_t, ip_recv_attr_t *);
3208 extern void     icmp_unreachable(mblk_t *, uint8_t, ip_recv_attr_t *);
3209 extern boolean_t ip_ipsec_policy_inherit(conn_t *, conn_t *, ip_recv_attr_t *);
3210 extern void     *ip_pullup(mblk_t *, ssize_t, ip_recv_attr_t *);
3211 extern void     ip_setl2src(mblk_t *, ip_recv_attr_t *, ill_t *);
3212 extern mblk_t   *ip_check_and_align_header(mblk_t *, uint_t, ip_recv_attr_t *);
3213 extern mblk_t   *ip_check_length(mblk_t *, uchar_t *, ssize_t, uint_t, uint_t,
3214     ip_recv_attr_t *);
3215 extern mblk_t   *ip_check_optlen(mblk_t *, ipha_t *, uint_t, uint_t,
3216     ip_recv_attr_t *);
3217 extern mblk_t   *ip_fix_dbref(mblk_t *, ip_recv_attr_t *);
3218 extern uint_t   ip_cksum(mblk_t *, int, uint32_t);
3219 extern int      ip_close(queue_t *, int);
3220 extern uint16_t ip_csum_hdr(ipha_t *);
3221 extern void     ip_forward_xmit_v4(nce_t *, ill_t *, mblk_t *, ipha_t *,
3222     ip_recv_attr_t *, uint32_t, uint32_t);
3223 extern boolean_t ip_forward_options(mblk_t *, ipha_t *, ill_t *,
3224     ip_recv_attr_t *);
3225 extern int      ip_fragment_v4(mblk_t *, nce_t *, iaflags_t, uint_t, uint32_t,
3226     uint32_t, zoneid_t, zoneid_t, pfirepostfrag_t postfragfn,
3227     uintptr_t *cookie);
3228 extern void     ip_proto_not_sup(mblk_t *, ip_recv_attr_t *);
3229 extern void     ip_ire_g_fini(void);
```

```
3230 extern void     ip_ire_g_init(void);
3231 extern void     ip_ire_fini(ip_stack_t *);
3232 extern void     ip_ire_init(ip_stack_t *);
3233 extern void     ip_mdata_to_mhi(ill_t *, mblk_t *, struct mac_header_info_s *);
3234 extern int      ip_openv4(queue_t *q, dev_t *devp, int flag, int sflag,
3235                 cred_t *credp);
3236 extern int      ip_openv6(queue_t *q, dev_t *devp, int flag, int sflag,
3237                 cred_t *credp);
3238 extern int      ip_reassemble(mblk_t *, ipf_t *, uint_t, boolean_t, ill_t *,
3239     size_t);
3240 extern void     ip_rput(queue_t *, mblk_t *);
3241 extern void     ip_input(ill_t *, ill_rx_ring_t *, mblk_t *,
3242     struct mac_header_info_s *);
3243 extern void     ip_input_v6(ill_t *, ill_rx_ring_t *, mblk_t *,
3244     struct mac_header_info_s *);
3245 extern mblk_t   *ip_input_common_v4(ill_t *, ill_rx_ring_t *, mblk_t *,
3246     struct mac_header_info_s *, squeue_t *, mblk_t **, uint_t *);
3247 extern mblk_t   *ip_input_common_v6(ill_t *, ill_rx_ring_t *, mblk_t *,
3248     struct mac_header_info_s *, squeue_t *, mblk_t **, uint_t *);
3249 extern void     ill_input_full_v4(mblk_t *, void *, void *,
3250     ip_recv_attr_t *, rtc_t *);
3251 extern void     ill_input_short_v4(mblk_t *, void *, void *,
3252     ip_recv_attr_t *, rtc_t *);
3253 extern void     ill_input_full_v6(mblk_t *, void *, void *,
3254     ip_recv_attr_t *, rtc_t *);
3255 extern void     ill_input_short_v6(mblk_t *, void *, void *,
3256     ip_recv_attr_t *, rtc_t *);
3257 extern ipaddr_t ip_input_options(ipha_t *, ipaddr_t, mblk_t *,
3258     ip_recv_attr_t *, int *);
3259 extern boolean_t ip_input_local_options(mblk_t *, ipha_t *, ip_recv_attr_t *);
3260 extern mblk_t   *ip_input_fragment(mblk_t *, ipha_t *, ip_recv_attr_t *);
3261 extern mblk_t   *ip_input_fragment_v6(mblk_t *, ip6_t *, ip6_frag_t *, uint_t,
3262     ip_recv_attr_t *);
3263 extern void     ip_input_post_ipsec(mblk_t *, ip_recv_attr_t *);
3264 extern void     ip_fanout_v4(mblk_t *, ipha_t *, ip_recv_attr_t *);
3265 extern void     ip_fanout_v6(mblk_t *, ip6_t *, ip_recv_attr_t *);
3266 extern void     ip_fanout_proto_conn(conn_t *, mblk_t *, ipha_t *, ip6_t *,
3267     ip_recv_attr_t *);
3268 extern void     ip_fanout_proto_v4(mblk_t *, ipha_t *, ip_recv_attr_t *);
3269 extern void     ip_fanout_send_icmp_v4(mblk_t *, uint_t, uint_t,
3270     ip_recv_attr_t *);
3271 extern void     ip_fanout_udp_conn(conn_t *, mblk_t *, ipha_t *, ip6_t *,
3272     ip_recv_attr_t *);
3273 extern void     ip_fanout_udp_multi_v4(mblk_t *, ipha_t *, uint16_t, uint16_t,
3274     ip_recv_attr_t *);
3275 extern mblk_t   *zero_spi_check(mblk_t *, ip_recv_attr_t *);
3276 extern void     ip_build_hdrs_v4(uchar_t *, uint_t, const ip_pkt_t *, uint8_t);
3277 extern int      ip_find_hdr_v4(ipha_t *, ip_pkt_t *, boolean_t);
3278 extern int      ip_total_hdrs_len_v4(const ip_pkt_t *);

3280 extern mblk_t   *ip_accept_tcp(ill_t *, ill_rx_ring_t *, squeue_t *,
3281     mblk_t *, mblk_t **, uint_t *cnt);
3282 extern void     ip_rput_dlpi(ill_t *, mblk_t *);
3283 extern void     ip_rput_notdata(ill_t *, mblk_t *);

3285 extern void     ip_mib2_add_ip_stats(mib2_ipIfStatsEntry_t *,
3286                 mib2_ipIfStatsEntry_t *);
3287 extern void     ip_mib2_add_icmp6_stats(mib2_ipv6IfIcmpEntry_t *,
3288                 mib2_ipv6IfIcmpEntry_t *);
3289 extern void     ip_rput_other(ipsq_t *, queue_t *, mblk_t *, void *);
3290 extern ire_t    *ip_check_multihome(void *, ire_t *, ill_t *);
3291 extern void     ip_send_potential_redirect_v4(mblk_t *, ipha_t *, ire_t *,
3292     ip_recv_attr_t *);
3293 extern int      ip_set_destination_v4(ipaddr_t *, ipaddr_t, ipaddr_t,
3294     ip_xmit_attr_t *, iulp_t *, uint32_t, uint_t);
3295 extern int      ip_set_destination_v6(in6_addr_t *, const in6_addr_t *,
```

```
3296          const in6_addr_t *, ip_xmit_attr_t *, iulp_t *, uint32_t, uint_t);

3298 extern int       ip_output_simple(mblk_t *, ip_xmit_attr_t *);
3299 extern int       ip_output_simple_v4(mblk_t *, ip_xmit_attr_t *);
3300 extern int       ip_output_simple_v6(mblk_t *, ip_xmit_attr_t *);
3301 extern int       ip_output_options(mblk_t *, ipha_t *, ip_xmit_attr_t *,
3302      ill_t *);
3303 extern void      ip_output_local_options(ipha_t *, ip_stack_t *);

3305 extern ip_xmit_attr_t *conn_get_ixa(conn_t *, boolean_t);
3306 extern ip_xmit_attr_t *conn_get_ixa_tryhard(conn_t *, boolean_t);
3307 extern ip_xmit_attr_t *conn_replace_ixa(conn_t *, ip_xmit_attr_t *);
3308 extern ip_xmit_attr_t *conn_get_ixa_exclusive(conn_t *);
3309 extern ip_xmit_attr_t *ip_xmit_attr_duplicate(ip_xmit_attr_t *);
3310 extern void      ip_xmit_attr_replace_tsl(ip_xmit_attr_t *, ts_label_t *);
3311 extern void      ip_xmit_attr_restore_tsl(ip_xmit_attr_t *, cred_t *);
3312 boolean_t        ip_recv_attr_replace_label(ip_recv_attr_t *, ts_label_t *);
3313 extern void      ixa_inactive(ip_xmit_attr_t *);
3314 extern void      ixa_refrele(ip_xmit_attr_t *);
3315 extern boolean_t ixa_check_drain_insert(conn_t *, ip_xmit_attr_t *);
3316 extern void      ixa_cleanup(ip_xmit_attr_t *);
3317 extern void      ira_cleanup(ip_recv_attr_t *, boolean_t);
3318 extern void      ixa_safe_copy(ip_xmit_attr_t *, ip_xmit_attr_t *);

3320 extern int       conn_ip_output(mblk_t *, ip_xmit_attr_t *);
3321 extern boolean_t ip_output_verify_local(ip_xmit_attr_t *);
3322 extern mblk_t    *ip_output_process_local(mblk_t *, ip_xmit_attr_t *, boolean_t,
3323      boolean_t, conn_t *);

3325 extern int       conn_opt_get(conn_opt_arg_t *, t_scalar_t, t_scalar_t,
3326      uchar_t *);
3327 extern int       conn_opt_set(conn_opt_arg_t *, t_scalar_t, t_scalar_t, uint_t,
3328      uchar_t *, boolean_t, cred_t *);
3329 extern boolean_t         conn_same_as_last_v4(conn_t *, sin_t *);
3330 extern boolean_t         conn_same_as_last_v6(conn_t *, sin6_t *);
3331 extern int       conn_update_label(const conn_t *, const ip_xmit_attr_t *,
3332      const in6_addr_t *, ip_pkt_t *);

3334 extern int       ip_opt_set_multicast_group(conn_t *, t_scalar_t,
3335      uchar_t *, boolean_t, boolean_t);
3336 extern int       ip_opt_set_multicast_sources(conn_t *, t_scalar_t,
3337      uchar_t *, boolean_t, boolean_t);
3338 extern int       conn_getsockname(conn_t *, struct sockaddr *, uint_t *);
3339 extern int       conn_getpeername(conn_t *, struct sockaddr *, uint_t *);

3341 extern int       conn_build_hdr_template(conn_t *, uint_t, uint_t,
3342      const in6_addr_t *, const in6_addr_t *, uint32_t);
3343 extern mblk_t    *conn_prepend_hdr(ip_xmit_attr_t *, const ip_pkt_t *,
3344      const in6_addr_t *, const in6_addr_t *, uint8_t, uint32_t, uint_t,
3345      mblk_t *, uint_t, uint_t, uint32_t *, int *);
3346 extern void      ip_attr_newdst(ip_xmit_attr_t *);
3347 extern void      ip_attr_nexthop(const ip_pkt_t *, const ip_xmit_attr_t *,
3348      const in6_addr_t *, in6_addr_t *);
3349 extern int       conn_connect(conn_t *, iulp_t *, uint32_t);
3350 extern int       ip_attr_connect(const conn_t *, ip_xmit_attr_t *,
3351      const in6_addr_t *, const in6_addr_t *, const in6_addr_t *, in_port_t,
3352      in6_addr_t *, iulp_t *, uint32_t);
3353 extern int       conn_inherit_parent(conn_t *, conn_t *);

3355 extern void      conn_ixa_cleanup(conn_t *connp, void *arg);

3357 extern boolean_t conn_wantpacket(conn_t *, ip_recv_attr_t *, ipha_t *);
3358 extern uint_t    ip_type_v4(ipaddr_t, ip_stack_t *);
3359 extern uint_t    ip_type_v6(const in6_addr_t *, ip_stack_t *);

3361 extern void      ip_wput_nondata(queue_t *, mblk_t *);
```

```
3362 extern void      ip_wsrv(queue_t *);
3363 extern char      *ip_nv_lookup(nv_t *, int);
3364 extern boolean_t ip_local_addr_ok_v6(const in6_addr_t *, const in6_addr_t *);
3365 extern boolean_t ip_remote_addr_ok_v6(const in6_addr_t *, const in6_addr_t *);
3366 extern ipaddr_t  ip_massage_options(ipha_t *, netstack_t *);
3367 extern ipaddr_t  ip_net_mask(ipaddr_t);
3368 extern void      arp_bringup_done(ill_t *, int);
3369 extern void      arp_replumb_done(ill_t *, int);

3371 extern struct qinit iprinitv6;

3373 extern void      ipmp_init(ip_stack_t *);
3374 extern void      ipmp_destroy(ip_stack_t *);
3375 extern ipmp_grp_t *ipmp_grp_create(const char *, phyint_t *);
3376 extern void      ipmp_grp_destroy(ipmp_grp_t *);
3377 extern void      ipmp_grp_info(const ipmp_grp_t *, lifgroupinfo_t *);
3378 extern int       ipmp_grp_rename(ipmp_grp_t *, const char *);
3379 extern ipmp_grp_t *ipmp_grp_lookup(const char *, ip_stack_t *);
3380 extern int       ipmp_grp_vet_phyint(ipmp_grp_t *, phyint_t *);
3381 extern ipmp_illgrp_t *ipmp_illgrp_create(ill_t *);
3382 extern void      ipmp_illgrp_destroy(ipmp_illgrp_t *);
3383 extern ill_t     *ipmp_illgrp_add_ipif(ipmp_illgrp_t *, ipif_t *);
3384 extern void      ipmp_illgrp_del_ipif(ipmp_illgrp_t *, ipif_t *);
3385 extern ill_t     *ipmp_illgrp_next_ill(ipmp_illgrp_t *);
3386 extern ill_t     *ipmp_illgrp_hold_next_ill(ipmp_illgrp_t *);
3387 extern ill_t     *ipmp_illgrp_hold_cast_ill(ipmp_illgrp_t *);
3388 extern ill_t     *ipmp_illgrp_ipmp_ill(ipmp_illgrp_t *);
3389 extern void      ipmp_illgrp_refresh_mtu(ipmp_illgrp_t *);
3390 extern ipmp_arpent_t *ipmp_illgrp_create_arpent(ipmp_illgrp_t *,
3391      boolean_t, ipaddr_t, uchar_t *, size_t, uint16_t);
3392 extern void      ipmp_illgrp_destroy_arpent(ipmp_illgrp_t *, ipmp_arpent_t *);
3393 extern ipmp_arpent_t *ipmp_illgrp_lookup_arpent(ipmp_illgrp_t *, ipaddr_t *);
3394 extern void      ipmp_illgrp_refresh_arpent(ipmp_illgrp_t *);
3395 extern void      ipmp_illgrp_mark_arpent(ipmp_illgrp_t *, ipmp_arpent_t *);
3396 extern ill_t     *ipmp_illgrp_find_ill(ipmp_illgrp_t *, uchar_t *, uint_t);
3397 extern void      ipmp_illgrp_link_grp(ipmp_illgrp_t *, ipmp_grp_t *);
3398 extern int       ipmp_illgrp_unlink_grp(ipmp_illgrp_t *);
3399 extern uint_t    ipmp_ill_get_ipmp_ifindex(const ill_t *);
3400 extern void      ipmp_ill_join_illgrp(ill_t *, ipmp_illgrp_t *);
3401 extern void      ipmp_ill_leave_illgrp(ill_t *);
3402 extern ill_t     *ipmp_ill_hold_ipmp_ill(ill_t *);
3403 extern ill_t     *ipmp_ill_hold_xmit_ill(ill_t *, boolean_t);
3404 extern boolean_t ipmp_ill_is_active(ill_t *);
3405 extern void      ipmp_ill_refresh_active(ill_t *);
3406 extern void      ipmp_phyint_join_grp(phyint_t *, ipmp_grp_t *);
3407 extern void      ipmp_phyint_leave_grp(phyint_t *);
3408 extern void      ipmp_phyint_refresh_active(phyint_t *);
3409 extern ill_t     *ipmp_ipif_bound_ill(const ipif_t *);
3410 extern ill_t     *ipmp_ipif_hold_bound_ill(const ipif_t *);
3411 extern boolean_t ipmp_ipif_is_dataaddr(const ipif_t *);
3412 extern boolean_t ipmp_ipif_is_stubaddr(const ipif_t *);
3413 extern boolean_t ipmp_packet_is_probe(mblk_t *, ill_t *);
3414 extern void      ipmp_ncec_delete_nce(ncec_t *);
3415 extern void      ipmp_ncec_refresh_nce(ncec_t *);

3417 extern void      conn_drain_insert(conn_t *, idl_tx_list_t *);
3418 extern void      conn_setqfull(conn_t *, boolean_t *);
3419 extern void      conn_clrqfull(conn_t *, boolean_t *);
3420 extern int       conn_ipsec_length(conn_t *);
3421 extern ipaddr_t  ip_get_dst(ipha_t *);
3422 extern uint_t    ip_get_pmtu(ip_xmit_attr_t *);
3423 extern uint_t    ip_get_base_mtu(ill_t *, ire_t *);
3424 extern mblk_t *ip_output_attach_policy(mblk_t *, ipha_t *, ip6_t *,
3425      const conn_t *, ip_xmit_attr_t *);
3426 extern int       ipsec_out_extra_length(ip_xmit_attr_t *);
3427 extern int       ipsec_out_process(mblk_t *, ip_xmit_attr_t *);
```

```
3428 extern int       ip_output_post_ipsec(mblk_t *, ip_xmit_attr_t *);
3429 extern void       ipsec_out_to_in(ip_xmit_attr_t *, ill_t *ill,
3430     ip_recv_attr_t *);

3432 extern void       ire_cleanup(ire_t *);
3433 extern void       ire_inactive(ire_t *);
3434 extern boolean_t irb_inactive(irb_t *);
3435 extern ire_t     *ire_unlink(irb_t *);

3437 #ifdef DEBUG
3438 extern  boolean_t th_trace_ref(const void *, ip_stack_t *);
3439 extern  void     th_trace_unref(const void *);
3440 extern  void     th_trace_cleanup(const void *, boolean_t);
3441 extern  void     ire_trace_ref(ire_t *);
3442 extern  void     ire_untrace_ref(ire_t *);
3443 #endif

3445 extern int       ip_srcid_insert(const in6_addr_t *, zoneid_t, ip_stack_t *);
3446 extern int       ip_srcid_remove(const in6_addr_t *, zoneid_t, ip_stack_t *);
3447 extern void      ip_srcid_find_id(uint_t, in6_addr_t *, zoneid_t, netstack_t *);
3448 extern uint_t    ip_srcid_find_addr(const in6_addr_t *, zoneid_t, netstack_t *);

3450 extern uint8_t ipoptp_next(ipoptp_t *);
3451 extern uint8_t ipoptp_first(ipoptp_t *, ipha_t *);
3452 extern int       ip_opt_get_user(conn_t *, uchar_t *);
3453 extern int       ipsec_req_from_conn(conn_t *, ipsec_req_t *, int);
3454 extern int       ip_snmp_get(queue_t *q, mblk_t *mctl, int level, boolean_t);
3455 extern int       ip_snmp_set(queue_t *q, int, int, uchar_t *, int);
3456 extern void      ip_process_ioctl(ipsq_t *, queue_t *, mblk_t *, void *);
3457 extern void      ip_quiesce_conn(conn_t *);
3458 extern  void     ip_reprocess_ioctl(ipsq_t *, queue_t *, mblk_t *, void *);
3459 extern void      ip_ioctl_finish(queue_t *, mblk_t *, int, int, ipsq_t *);

3461 extern boolean_t ip_cmpbuf(const void *, uint_t, boolean_t, const void *,
3462     uint_t);
3463 extern boolean_t ip_allocbuf(void **, uint_t *, boolean_t, const void *,
3464     uint_t);
3465 extern void      ip_savebuf(void **, uint_t *, boolean_t, const void *, uint_t);

3467 extern boolean_t       ipsq_pending_mp_cleanup(ill_t *, conn_t *);
3468 extern void      conn_ioctl_cleanup(conn_t *);

3470 extern void      ip_unbind(conn_t *);

3472 extern void tnet_init(void);
3473 extern void tnet_fini(void);

3475 /*
3476  * Hook functions to enable cluster networking
3477  * On non-clustered systems these vectors must always be NULL.
3478  */
3479 extern int (*cl_inet_isclusterwide)(netstackid_t stack_id, uint8_t protocol,
3480     sa_family_t addr_family, uint8_t *laddrp, void *args);
3481 extern uint32_t (*cl_inet_ipident)(netstackid_t stack_id, uint8_t protocol,
3482     sa_family_t addr_family, uint8_t *laddrp, uint8_t *faddrp,
3483     void *args);
3484 extern int (*cl_inet_connect2)(netstackid_t stack_id, uint8_t protocol,
3485     boolean_t is_outgoing, sa_family_t addr_family, uint8_t *laddrp,
3486     in_port_t lport, uint8_t *faddrp, in_port_t fport, void *args);
3487 extern void (*cl_inet_getspi)(netstackid_t, uint8_t, uint8_t *, size_t,
3488     void *);
3489 extern void (*cl_inet_getspi)(netstackid_t stack_id, uint8_t protocol,
3490     uint8_t *ptr, size_t len, void *args);
3491 extern int (*cl_inet_checkspi)(netstackid_t stack_id, uint8_t protocol,
3492     uint32_t spi, void *args);
3493 extern void (*cl_inet_deletespi)(netstackid_t stack_id, uint8_t protocol,
```

```
3494     uint32_t spi, void *args);
3495 extern void (*cl_inet_idlesa)(netstackid_t, uint8_t, uint32_t,
3496     sa_family_t, in6_addr_t, in6_addr_t, void *);


3499 /* Hooks for CGTP (multirt routes) filtering module */
3500 #define CGTP_FILTER_REV_1        1
3501 #define CGTP_FILTER_REV_2        2
3502 #define CGTP_FILTER_REV_3        3
3503 #define CGTP_FILTER_REV          CGTP_FILTER_REV_3

3505 /* cfo_filter and cfo_filter_v6 hooks return values */
3506 #define CGTP_IP_PKT_NOT_CGTP     0
3507 #define CGTP_IP_PKT_PREMIUM      1
3508 #define CGTP_IP_PKT_DUPLICATE    2

3510 /* Version 3 of the filter interface */
3511 typedef struct cgtp_filter_ops {
3512         int     cfo_filter_rev;                 /* CGTP_FILTER_REV_3 */
3513         int     (*cfo_change_state)(netstackid_t, int);
3514         int     (*cfo_add_dest_v4)(netstackid_t, ipaddr_t, ipaddr_t,
3515                     ipaddr_t, ipaddr_t);
3516         int     (*cfo_del_dest_v4)(netstackid_t, ipaddr_t, ipaddr_t);
3517         int     (*cfo_add_dest_v6)(netstackid_t, in6_addr_t *, in6_addr_t *,
3518                     in6_addr_t *, in6_addr_t *);
3519         int     (*cfo_del_dest_v6)(netstackid_t, in6_addr_t *, in6_addr_t *);
3520         int     (*cfo_filter)(netstackid_t, uint_t, mblk_t *);
3521         int     (*cfo_filter_v6)(netstackid_t, uint_t, ip6_t *,
3522                     ip6_frag_t *);
3523 } cgtp_filter_ops_t;

3525 #define CGTP_MCAST_SUCCESS       1

3527 /*
3528  * The separate CGTP module needs this global symbol so that it
3529  * can check the version and determine whether to use the old or the new
3530  * version of the filtering interface.
3531  */
3532 extern int       ip_cgtp_filter_rev;

3534 extern int       ip_cgtp_filter_supported(void);
3535 extern int       ip_cgtp_filter_register(netstackid_t, cgtp_filter_ops_t *);
3536 extern int       ip_cgtp_filter_unregister(netstackid_t);
3537 extern int       ip_cgtp_filter_is_registered(netstackid_t);

3539 /*
3540  * rr_ring_state cycles in the order shown below from RR_FREE through
3541  * RR_FREE_IN_PROG and  back to RR_FREE.
3542  */
3543 typedef enum {
3544         RR_FREE,                        /* Free slot */
3545         RR_SQUEUE_UNBOUND,              /* Ring's squeue is unbound */
3546         RR_SQUEUE_BIND_INPROG,          /* Ring's squeue bind in progress */
3547         RR_SQUEUE_BOUND,                /* Ring's squeue bound to cpu */
3548         RR_FREE_INPROG                  /* Ring is being freed */
3549 } ip_ring_state_t;

3551 #define ILL_MAX_RINGS            256    /* Max num of rx rings we can manage */
3552 #define ILL_POLLING              0x01   /* Polling in use */

3554 /*
3555  * These functions pointer types are exported by the mac/dls layer.
3556  * we need to duplicate the definitions here because we cannot
3557  * include mac/dls header files here.
3558  */
3559 typedef boolean_t               (*ip_mac_intr_disable_t)(void *);
```

```
3560 typedef void                      (*ip_mac_intr_enable_t)(void *);
3561 typedef ip_mac_tx_cookie_t        (*ip_dld_tx_t)(void *, mblk_t *,
3562     uint64_t, uint16_t);
3563 typedef void                      (*ip_flow_enable_t)(void *, ip_mac_tx_cookie_t);
3564 typedef void                      *(*ip_dld_callb_t)(void *,
3565     ip_flow_enable_t, void *);
3566 typedef boolean_t                 (*ip_dld_fctl_t)(void *, ip_mac_tx_cookie_t);
3567 typedef int                       (*ip_capab_func_t)(void *, uint_t,
3568     void *, uint_t);

3570 /*
3571  * POLLING README
3572  * sq_get_pkts() is called to pick packets from softring in poll mode. It
3573  * calls rr_rx to get the chain and process it with rr_ip_accept.
3574  * rr_rx = mac_soft_ring_poll() to pick packets
3575  * rr_ip_accept = ip_accept_tcp() to process packets
3576  */

3578 /*
3579  * XXX: With protocol, service specific squeues, they will have
3580  * specific acceptor functions.
3581  */
3582 typedef mblk_t *(*ip_mac_rx_t)(void *, size_t);
3583 typedef mblk_t *(*ip_accept_t)(ill_t *, ill_rx_ring_t *,
3584     squeue_t *, mblk_t *, mblk_t **, uint_t *);

3586 /*
3587  * rr_intr_enable, rr_intr_disable, rr_rx_handle, rr_rx:
3588  * May be accessed while in the squeue AND after checking that SQS_POLL_CAPAB
3589  * is set.
3590  *
3591  * rr_ring_state: Protected by ill_lock.
3592  */
3593 struct ill_rx_ring {
3594         ip_mac_intr_disable_t   rr_intr_disable; /* Interrupt disabling func */
3595         ip_mac_intr_enable_t    rr_intr_enable; /* Interrupt enabling func */
3596         void                    *rr_intr_handle; /* Handle interrupt funcs */
3597         ip_mac_rx_t             rr_rx;          /* Driver receive function */
3598         ip_accept_t             rr_ip_accept;   /* IP accept function */
3599         void                    *rr_rx_handle;  /* Handle for Rx ring */
3600         squeue_t                *rr_sqp; /* Squeue the ring is bound to */
3601         ill_t                   *rr_ill;        /* back pointer to ill */
3602         ip_ring_state_t         rr_ring_state;  /* State of this ring */
3603 };

3605 /*
3606  * IP - DLD direct function call capability
3607  * Suffixes, df - dld function, dh - dld handle,
3608  * cf - client (IP) function, ch - client handle
3609  */
3610 typedef struct ill_dld_direct_s {                 /* DLD provided driver Tx */
3611         ip_dld_tx_t             idd_tx_df;      /* str_mdata_fastpath_put */
3612         void                    *idd_tx_dh;     /* dld_str_t *dsp */
3613         ip_dld_callb_t          idd_tx_cb_df;   /* mac_tx_srs_notify */
3614         void                    *idd_tx_cb_dh;  /* mac_client_handle_t *mch */
3615         ip_dld_fctl_t           idd_tx_fctl_df; /* mac_tx_is_flow_blocked */
3616         void                    *idd_tx_fctl_dh;        /* mac_client_handle */
3617 } ill_dld_direct_t;

3619 /* IP - DLD polling capability */
3620 typedef struct ill_dld_poll_s {
3621         ill_rx_ring_t           idp_ring_tbl[ILL_MAX_RINGS];
3622 } ill_dld_poll_t;

3624 /* Describes ill->ill_dld_capab */
3625 struct ill_dld_capab_s {
```

```
3626         ip_capab_func_t         idc_capab_df;   /* dld_capab_func */
3627         void                    *idc_capab_dh;  /* dld_str_t *dsp */
3628         ill_dld_direct_t        idc_direct;
3629         ill_dld_poll_t          idc_poll;
3630 };

3632 /*
3633  * IP squeues exports
3634  */
3635 extern boolean_t        ip_squeue_fanout;

3637 #define IP_SQUEUE_GET(hint) ip_squeue_random(hint)

3639 extern void ip_squeue_init(void (*)(squeue_t *));
3640 extern squeue_t *ip_squeue_random(uint_t);
3641 extern squeue_t *ip_squeue_get(ill_rx_ring_t *);
3642 extern squeue_t *ip_squeue_getfree(pri_t);
3643 extern int ip_squeue_cpu_move(squeue_t *, processorid_t);
3644 extern void *ip_squeue_add_ring(ill_t *, void *);
3645 extern void ip_squeue_bind_ring(ill_t *, ill_rx_ring_t *, processorid_t);
3646 extern void ip_squeue_clean_ring(ill_t *, ill_rx_ring_t *);
3647 extern void ip_squeue_quiesce_ring(ill_t *, ill_rx_ring_t *);
3648 extern void ip_squeue_restart_ring(ill_t *, ill_rx_ring_t *);
3649 extern void ip_squeue_clean_all(ill_t *);
3650 extern boolean_t        ip_source_routed(ipha_t *, ip_stack_t *);

3652 extern void tcp_wput(queue_t *, mblk_t *);

3654 extern int      ip_fill_mtuinfo(conn_t *, ip_xmit_attr_t *,
3655     struct ip6_mtuinfo *);
3656 extern hook_t *ipobs_register_hook(netstack_t *, pfv_t);
3657 extern void ipobs_unregister_hook(netstack_t *, hook_t *);
3658 extern void ipobs_hook(mblk_t *, int, zoneid_t, zoneid_t, const ill_t *,
3659     ip_stack_t *);
3660 typedef void    (*ipsq_func_t)(ipsq_t *, queue_t *, mblk_t *, void *);

3662 extern void     dce_g_init(void);
3663 extern void     dce_g_destroy(void);
3664 extern void     dce_stack_init(ip_stack_t *);
3665 extern void     dce_stack_destroy(ip_stack_t *);
3666 extern void     dce_cleanup(uint_t, ip_stack_t *);
3667 extern dce_t    *dce_get_default(ip_stack_t *);
3668 extern dce_t    *dce_lookup_pkt(mblk_t *, ip_xmit_attr_t *, uint_t *);
3669 extern dce_t    *dce_lookup_v4(ipaddr_t, ip_stack_t *, uint_t *);
3670 extern dce_t    *dce_lookup_v6(const in6_addr_t *, uint_t, ip_stack_t *,
3671     uint_t *);
3672 extern dce_t    *dce_lookup_and_add_v4(ipaddr_t, ip_stack_t *);
3673 extern dce_t    *dce_lookup_and_add_v6(const in6_addr_t *, uint_t,
3674     ip_stack_t *);
3675 extern int      dce_update_uinfo_v4(ipaddr_t, iulp_t *, ip_stack_t *);
3676 extern int      dce_update_uinfo_v6(const in6_addr_t *, uint_t, iulp_t *,
3677     ip_stack_t *);
3678 extern int      dce_update_uinfo(const in6_addr_t *, uint_t, iulp_t *,
3679     ip_stack_t *);
3680 extern void     dce_increment_generation(dce_t *);
3681 extern void     dce_increment_all_generations(boolean_t, ip_stack_t *);
3682 extern void     dce_refrele(dce_t *);
3683 extern void     dce_refhold(dce_t *);
3684 extern void     dce_refrele_notr(dce_t *);
3685 extern void     dce_refhold_notr(dce_t *);
3686 mblk_t          *ip_snmp_get_mib2_ip_dce(queue_t *, mblk_t *, ip_stack_t *ipst);

3688 extern ip_laddr_t ip_laddr_verify_v4(ipaddr_t, zoneid_t,
3689     ip_stack_t *, boolean_t);
3690 extern ip_laddr_t ip_laddr_verify_v6(const in6_addr_t *, zoneid_t,
3691     ip_stack_t *, boolean_t, uint_t);
```

```
3692 extern int      ip_laddr_fanout_insert(conn_t *);

3694 extern boolean_t ip_verify_src(mblk_t *, ip_xmit_attr_t *, uint_t *);
3695 extern int      ip_verify_ire(mblk_t *, ip_xmit_attr_t *);

3697 extern mblk_t   *ip_xmit_attr_to_mblk(ip_xmit_attr_t *);
3698 extern boolean_t ip_xmit_attr_from_mblk(mblk_t *, ip_xmit_attr_t *);
3699 extern mblk_t   *ip_xmit_attr_free_mblk(mblk_t *);
3700 extern mblk_t   *ip_recv_attr_to_mblk(ip_recv_attr_t *);
3701 extern boolean_t ip_recv_attr_from_mblk(mblk_t *, ip_recv_attr_t *);
3702 extern mblk_t   *ip_recv_attr_free_mblk(mblk_t *);
3703 extern boolean_t ip_recv_attr_is_mblk(mblk_t *);

3705 extern char     *inet_ntop(int, const void *, char *, int);
3706 extern int      _inet_pton(int, char *, void *);
3707 #define inet_pton(x, y, z)     _inet_pton(x, y, z)

3709 #endif /* ! codereview */
3710 /*
3711  * Squeue tags. Tags only need to be unique when the callback function is the
3712  * same to distinguish between different calls, but we use unique tags for
3713  * convenience anyway.
3714  */
3715 #define SQTAG_IP_INPUT              1
3716 #define SQTAG_TCP_INPUT_ICMP_ERR    2
3717 #define SQTAG_TCP6_INPUT_ICMP_ERR   3
3718 #define SQTAG_IP_TCP_INPUT          4
3719 #define SQTAG_IP6_TCP_INPUT         5
3720 #define SQTAG_IP_TCP_CLOSE          6
3721 #define SQTAG_TCP_OUTPUT            7
3722 #define SQTAG_TCP_TIMER             8
3723 #define SQTAG_TCP_TIMEWAIT          9
3724 #define SQTAG_TCP_ACCEPT_FINISH     10
3725 #define SQTAG_TCP_ACCEPT_FINISH_Q0  11
3726 #define SQTAG_TCP_ACCEPT_PENDING    12
3727 #define SQTAG_TCP_LISTEN_DISCON     13
3728 #define SQTAG_TCP_CONN_REQ_1        14
3729 #define SQTAG_TCP_EAGER_BLOWOFF     15
3730 #define SQTAG_TCP_EAGER_CLEANUP     16
3731 #define SQTAG_TCP_EAGER_CLEANUP_Q0  17
3732 #define SQTAG_TCP_CONN_IND          18
3733 #define SQTAG_TCP_RSRV              19
3734 #define SQTAG_TCP_ABORT_BUCKET      20
3735 #define SQTAG_TCP_REINPUT          21
3736 #define SQTAG_TCP_REINPUT_EAGER    22
3737 #define SQTAG_TCP_INPUT_MCTL       23
3738 #define SQTAG_TCP_RPUTOTHER        24
3739 #define SQTAG_IP_PROTO_AGAIN       25
3740 #define SQTAG_IP_FANOUT_TCP        26
3741 #define SQTAG_IPSQ_CLEAN_RING      27
3742 #define SQTAG_TCP_WPUT_OTHER       28
3743 #define SQTAG_TCP_CONN_REQ_UNBOUND 29
3744 #define SQTAG_TCP_SEND_PENDING     30
3745 #define SQTAG_BIND_RETRY           31
3746 #define SQTAG_UDP_FANOUT           32
3747 #define SQTAG_UDP_INPUT            33
3748 #define SQTAG_UDP_WPUT             34
3749 #define SQTAG_UDP_OUTPUT           35
3750 #define SQTAG_TCP_KSSL_INPUT       36
3751 #define SQTAG_TCP_DROP_Q0          37
3752 #define SQTAG_TCP_CONN_REQ_2       38
3753 #define SQTAG_IP_INPUT_RX_RING     39
3754 #define SQTAG_SQUEUE_CHANGE        40
3755 #define SQTAG_CONNECT_FINISH       41
3756 #define SQTAG_SYNCHRONOUS_OP       42
3757 #define SQTAG_TCP_SHUTDOWN_OUTPUT  43
```

```
3758 #define SQTAG_TCP_IXA_CLEANUP           44
3759 #define SQTAG_TCP_SEND_SYNACK           45

3761 extern sin_t    sin_null;       /* Zero address for quick clears */
3762 extern sin6_t   sin6_null;      /* Zero address for quick clears */

3764 #endif  /* _KERNEL */

3766 #ifdef  __cplusplus
3767 }
3768 #endif

3770 #endif  /* _INET_IP_H */
```

```
   **********************************************************
      10608 Wed Sep 26 12:51:44 2012
   new/usr/src/uts/common/inet/ip/inet_ntop.c
   inet_pton
   **********************************************************
     1 /*
     2  * CDDL HEADER START
     3  *
     4  * The contents of this file are subject to the terms of the
     5  * Common Development and Distribution License, Version 1.0 only
     6  * (the "License").  You may not use this file except in compliance
     7  * with the License.
     8  *
     9  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
    10  * or http://www.opensolaris.org/os/licensing.
    11  * See the License for the specific language governing permissions
    12  * and limitations under the License.
    13  *
    14  * When distributing Covered Code, include this CDDL HEADER in each
    15  * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
    16  * If applicable, add the following below this CDDL HEADER, with the
    17  * fields enclosed by brackets "[]" replaced with your own identifying
    18  * information: Portions Copyright [yyyy] [name of copyright owner]
    19  *
    20  * CDDL HEADER END
    21  */
    22 /*
    23  * Copyright 2004 Sun Microsystems, Inc.  All rights reserved.
    24  * Use is subject to license terms.
    25  */
    26 /*
    27  * Copyright 2012 Nexenta Systems, Inc. All rights reserved.
    28  */

    27 #pragma ident   "%Z%%M% %I%     %E% SMI"

    30 #include <sys/types.h>
    31 #include <sys/cmn_err.h>
    32 #include <sys/systm.h>
    33 #include <sys/socket.h>
    34 #include <sys/sunddi.h>
    35 #include <netinet/in.h>
    36 #include <inet/led.h>

    38 static void     convert2ascii(char *, const in6_addr_t *);
    39 static char     *strchr_w(const char *, int);
    40 static int      str2inet_addr(char *, ipaddr_t *);

    42 /*
    43  * inet_ntop -- Convert an IPv4 or IPv6 address in binary form into
    44  * printable form, and return a pointer to that string. Caller should
    45  * provide a buffer of correct length to store string into.
    46  * Note: this routine is kernel version of inet_ntop. It has similar
    47  * format as inet_ntop() defined in rfc2553. But it does not do
    48  * error handling operations exactly as rfc2553 defines. This function
    49  * is used by kernel inet directory routines only for debugging.
    50  * This inet_ntop() function, does not return NULL if third argument
    51  * is NULL. The reason is simple that we don't want kernel to panic
    52  * as the output of this function is directly fed to ip<n>dbg macro.
    53  * Instead it uses a local buffer for destination address for
    54  * those calls which purposely pass NULL ptr for the destination
    55  * buffer. This function is thread-safe when the caller passes a non-
    56  * null buffer with the third argument.
    57  */
    58 /* ARGSUSED */
    59 char *
```

```
    60 inet_ntop(int af, const void *addr, char *buf, int addrlen)
    61 {
    62         static char local_buf[INET6_ADDRSTRLEN];
    63         static char *err_buf1 = "<badaddr>";
    64         static char *err_buf2 = "<badfamily>";
    65         in6_addr_t      *v6addr;
    66         uchar_t         *v4addr;
    67         char            *caddr;

    69         /*
    70          * We don't allow thread unsafe inet_ntop calls, they
    71          * must pass a non-null buffer pointer. For DEBUG mode
    72          * we use the ASSERT() and for non-debug kernel it will
    73          * silently allow it for now. Someday we should remove
    74          * the static buffer from this function.
    75          */

    77         ASSERT(buf != NULL);
    78         if (buf == NULL)
    79                 buf = local_buf;
    80         buf[0] = '\0';

    82         /* Let user know politely not to send NULL or unaligned addr */
    83         if (addr == NULL || !(OK_32PTR(addr))) {
    84 #ifdef DEBUG
    85                 cmn_err(CE_WARN, "inet_ntop: addr is <null> or unaligned");
    86 #endif
    87                 return (err_buf1);
    88         }

    91 #define UC(b)   (((int)b) & 0xff)
    92         switch (af) {
    93         case AF_INET:
    94                 ASSERT(addrlen >= INET_ADDRSTRLEN);
    95                 v4addr = (uchar_t *)addr;
    96                 (void) sprintf(buf, "%03d.%03d.%03d.%03d",
    97                     UC(v4addr[0]), UC(v4addr[1]), UC(v4addr[2]), UC(v4addr[3]));
    98                 return (buf);

   100         case AF_INET6:
   101                 ASSERT(addrlen >= INET6_ADDRSTRLEN);
   102                 v6addr = (in6_addr_t *)addr;
   103                 if (IN6_IS_ADDR_V4MAPPED(v6addr)) {
   104                         caddr = (char *)addr;
   105                         (void) sprintf(buf, "::ffff:%d.%d.%d.%d",
   106                             UC(caddr[12]), UC(caddr[13]),
   107                             UC(caddr[14]), UC(caddr[15]));
   108                 } else if (IN6_IS_ADDR_V4COMPAT(v6addr)) {
   109                         caddr = (char *)addr;
   110                         (void) sprintf(buf, "::%d.%d.%d.%d",
   111                             UC(caddr[12]), UC(caddr[13]), UC(caddr[14]),
   112                             UC(caddr[15]));
   113                 } else if (IN6_IS_ADDR_UNSPECIFIED(v6addr)) {
   114                         (void) sprintf(buf, "::");
   115                 } else {
   116                         convert2ascii(buf, v6addr);
   117                 }
   118                 return (buf);

   120         default:
   121                 return (err_buf2);
   122         }
   123 #undef UC
   124 }
   _____unchanged_portion_omitted_
```

```
 231 static int
 232 str2inet_addr(char *cp, ipaddr_t *addrp)
 233 {
 234         char *end;
 235         long byte;
 236         int i;
 237         uint8_t *addr = (uint8_t *)addrp;

 239         *addrp = 0;
 236         ipaddr_t addr = 0;

 241         for (i = 0; i < 4; i++) {
 242                 if (ddi_strtol(cp, &end, 10, &byte) != 0 || byte < 0 ||
 243                     byte > 255) {
 244                         return (0);
 245                 }
 246                 addr[i] = (uint8_t)byte;
 243                 addr = (addr << 8) | (uint8_t)byte;
 247                 if (i < 3) {
 248                         if (*end != '.') {
 249                                 return (0);
 250                         } else {
 251                                 cp = end + 1;
 252                         }
 253                 } else {
 254                         cp = end;
 255                 }
 256         }

 254         *addrp = addr;
 258         return (1);
 259 }

 261 /*
 262  * inet_pton: This function takes string format IPv4 or IPv6 address and
 263  * converts it to binary form. The format of this function corresponds to
 264  * inet_pton() in the socket library.
 265  * It returns 0 for invalid IPv4 and IPv6 address
 266  *            1 when successfully converts ascii to binary
 267  *            -1 when af is not AF_INET or AF_INET6
 268  */
 269 int
 270 m_inet_pton(int af, char *inp, void *outp, int revert)
 267 inet_pton(int af, char *inp, void *outp)
 271 {
 272         int i;
 273         long byte;
 274         char *end;

 276         switch (af) {
 277         case AF_INET:
 278                 if (str2inet_addr(inp, (ipaddr_t *)outp)) {
 279                         if (! revert)
 280                                 *(uint32_t *)outp = ntohl(*(uint32_t *)outp);
 281                         return (1);
 282                 } else {
 283                         return (0);
 284                 }
 275                 return (str2inet_addr(inp, (ipaddr_t *)outp));
 285         case AF_INET6: {
 286                 union v6buf_u {
 287                         uint16_t v6words_u[8];
 288                         in6_addr_t v6addr_u;
 289                 } v6buf, *v6outp;
 290                 uint16_t        *dbl_col = NULL;
```

```
 291                 char lastbyte = NULL;

 293                 v6outp = (union v6buf_u *)outp;

 295                 if (strchr_w(inp, '.') != NULL) {
 296                         /* v4 mapped or v4 compatable */
 297                         if (strncmp(inp, "::ffff:", 7) == 0) {
 298                                 ipaddr_t ipv4_all_zeroes = 0;
 299                                 /* mapped - first init prefix and then fill */
 300                                 IN6_IPADDR_TO_V4MAPPED(ipv4_all_zeroes,
 301                                     &v6outp->v6addr_u);
 302                                 return (str2inet_addr(inp + 7,
 303                                     &(v6outp->v6addr_u.s6_addr32[3])));
 304                         } else if (strncmp(inp, "::", 2) == 0) {
 305                                 /* v4 compatable - prefix all zeroes */
 306                                 bzero(&v6outp->v6addr_u, sizeof (in6_addr_t));
 307                                 return (str2inet_addr(inp + 2,
 308                                     &(v6outp->v6addr_u.s6_addr32[3])));
 309                         }
 310                         return (0);
 311                 }
 312                 for (i = 0; i < 8; i++) {
 313                         int error;
 314                         /*
 315                          * if ddi_strtol() fails it could be because
 316                          * the string is "::".  That is valid and
 317                          * checked for below so just set the value to
 318                          * 0 and continue.
 319                          */
 320                         if ((error = ddi_strtol(inp, &end, 16, &byte)) != 0) {
 321                                 if (error == ERANGE)
 322                                         return (0);
 323                                 byte = 0;
 324                         }
 325                         if (byte < 0 || byte > 0x0ffff) {
 326                                 return (0);
 327                         }
 328                         if (revert) {
 329                                 v6buf.v6words_u[i] = htons((uint16_t)byte);
 330                         } else {
 331 #endif /* ! codereview */
 332                                 v6buf.v6words_u[i] = (uint16_t)byte;
 333                         }
 334 #endif /* ! codereview */
 335                         if (*end == NULL || i == 7) {
 336                                 inp = end;
 337                                 break;
 338                         }
 339                         if (inp == end) {       /* not a number must be */
 340                                 if (*inp == ':' &&
 341                                     ((i == 0 && *(inp + 1) == ':') ||
 342                                     lastbyte == ':')) {
 343                                         if (dbl_col) {
 344                                                 return (0);
 345                                         }
 346                                         if (byte != 0)
 347                                                 i++;
 348                                         dbl_col = &v6buf.v6words_u[i];
 349                                         if (i == 0)
 350                                                 inp++;
 351                                 } else if (*inp == NULL || *inp == ' ' ||
 352                                     *inp == '\t') {
 353                                         break;
 354                                 } else {
 355                                         return (0);
 356                                 }
```

```
357                          } else {
358                                  inp = end;
359                          }
360                          if (*inp != ':') {
361                                  return (0);
362                          }
363                          inp++;
364                          if (*inp == NULL || *inp == ' ' || *inp == '\t') {
365                                  break;
366                          }
367                          lastbyte = *inp;
368                  }
369                  if (*inp != NULL && *inp != ' ' && *inp != '\t') {
370                          return (0);
371                  }
372                  /*
373                   * v6words now contains the bytes we could translate
374                   * dbl_col points to the word (should be 0) where
375                   * a double colon was found
376                   */
377                  if (i == 7) {
378                          v6outp->v6addr_u = v6buf.v6addr_u;
379                  } else {
380                          int rem;
381                          int word;
382                          int next;
383                          if (dbl_col == NULL) {
384                                  return (0);
385                          }
386                          bzero(&v6outp->v6addr_u, sizeof (in6_addr_t));
387                          rem = dbl_col - &v6buf.v6words_u[0];
388                          for (next = 0; next < rem; next++) {
389                                  v6outp->v6words_u[next] = v6buf.v6words_u[next];
390                          }
391                          next++; /* skip dbl_col 0 */
392                          rem = i - rem;
393                          word = 8 - rem;
394                          while (rem > 0) {
395                                  v6outp->v6words_u[word] = v6buf.v6words_u[next];
396                                  word++;
397                                  rem--;
398                                  next++;
399                          }
400                  }
401                  return (1);      /* Success */
402          }
403          }               /* switch */
404          return (-1);     /* return -1 for default case */
405 }


408 int
409 _inet_pton(int af, char *inp, void *outp)
410 {
411          return (m_inet_pton(af, inp, outp, 1));
412 }

414 /*
415  * We need this inet_pton to preserve compatibility with old closed binaries.
416  * Earlier, inet_pton returned address in hardware native order,
417  * not in network one. (See http://www.illumos.org/issue/3105).
418  * Having fixed that, we still need to support binaries, that use bad inet_pton
419  * and reverse returned address manually. All new inet_pton calls will be
420  * redirected to _inet_pton with #define in the header file.
421  */
422 int
```

```
423 inet_pton(int af, char *inp, void *outp)
424 {
425          return (m_inet_pton(af, inp, outp, 0));
426 }
427 #endif /* ! codereview */
```

```
**********************************************************
   10051 Wed Sep 26 12:51:46 2012
new/usr/src/uts/common/inet/ip6.h
inet_pton
**********************************************************
   1 /*
   2  * CDDL HEADER START
   3  *
   4  * The contents of this file are subject to the terms of the
   5  * Common Development and Distribution License (the "License").
   6  * You may not use this file except in compliance with the License.
   7  *
   8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
   9  * or http://www.opensolaris.org/os/licensing.
  10  * See the License for the specific language governing permissions
  11  * and limitations under the License.
  12  *
  13  * When distributing Covered Code, include this CDDL HEADER in each
  14  * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
  15  * If applicable, add the following below this CDDL HEADER, with the
  16  * fields enclosed by brackets "[]" replaced with your own identifying
  17  * information: Portions Copyright [yyyy] [name of copyright owner]
  18  *
  19  * CDDL HEADER END
  20  */
  21 /*
  22  * Copyright 2009 Sun Microsystems, Inc.  All rights reserved.
  23  * Use is subject to license terms.
  24  */
  25 /*
  26  * Copyright 2012 Nexenta Systems, Inc. All rights reserved.
  27  */
  28 #endif /* ! codereview */

  30 #ifndef _INET_IP6_H
  31 #define _INET_IP6_H

  33 #ifdef  __cplusplus
  34 extern "C" {
  35 #endif

  37 #include <sys/isa_defs.h>

  39 #ifdef  _KERNEL
  40 /* icmp6_t is used in the prototype of icmp_inbound_error_fanout_v6() */
  41 #include <netinet/icmp6.h>
  42 #endif  /* _KERNEL */

  44 /* version number for IPv6 - hard to get this one wrong! */
  45 #define IPV6_VERSION            6

  47 #define IPV6_HDR_LEN            40

  49 #define IPV6_ADDR_LEN           16

  51 /*
  52  * IPv6 address scopes.  The values of these enums also match the scope
  53  * field of multicast addresses.
  54  */
  55 typedef enum {
  56         IP6_SCOPE_INTFLOCAL = 1,        /* Multicast addresses only */
  57         IP6_SCOPE_LINKLOCAL,
  58         IP6_SCOPE_SUBNETLOCAL,          /* Multicast addresses only */
  59         IP6_SCOPE_ADMINLOCAL,           /* Multicast addresses only */
  60         IP6_SCOPE_SITELOCAL,
  61         IP6_SCOPE_GLOBAL
```

```
  62 } in6addr_scope_t;

  64 /* From RFC 3542 - setting for IPV6_USE_MIN_MTU socket option */
  65 #define IPV6_USE_MIN_MTU_MULTICAST      -1      /* Default */
  66 #define IPV6_USE_MIN_MTU_NEVER          0
  67 #define IPV6_USE_MIN_MTU_ALWAYS         1

  69 #ifdef  _KERNEL

  71 /* Extract the scope from a multicast address */
  72 #ifdef _BIG_ENDIAN
  73 #define IN6_ADDR_MC_SCOPE(addr) \
  74         (((addr)->s6_addr32[0] & 0x000f0000) >> 16)
  75 #else
  76 #define IN6_ADDR_MC_SCOPE(addr) \
  77         (((addr)->s6_addr32[0] & 0x00000f00) >> 8)
  78 #endif

  80 /* Default IPv4 TTL for IPv6-in-IPv4 encapsulated packets */
  81 #define IPV6_DEFAULT_HOPS       60      /* XXX What should it be? */

  83 /* Max IPv6 TTL */
  84 #define IPV6_MAX_HOPS   255

  86 /* Minimum IPv6 MTU from rfc2460 */
  87 #define IPV6_MIN_MTU            1280

  89 /* EUI-64 based token length */
  90 #define IPV6_TOKEN_LEN          64

  92 /* Length of an advertised IPv6 prefix */
  93 #define IPV6_PREFIX_LEN         64

  95 /* Default and maximum tunnel encapsulation limits.  See RFC 2473. */
  96 #define IPV6_DEFAULT_ENCAPLIMIT 4
  97 #define IPV6_MAX_ENCAPLIMIT     255

  99 /*
 100  * Minimum and maximum extension header lengths for IPv6.  The 8-bit
 101  * length field of each extension header (see rfc2460) specifies the
 102  * number of 8 octet units of data in the header not including the
 103  * first 8 octets.  A value of 0 would indicate 8 bytes (0 * 8 + 8),
 104  * and 255 would indicate 2048 bytes (255 * 8 + 8).
 105  */
 106 #define MIN_EHDR_LEN            8
 107 #define MAX_EHDR_LEN            2048

 109 #ifdef _BIG_ENDIAN
 110 #define IPV6_DEFAULT_VERS_AND_FLOW      0x60000000
 111 #define IPV6_VERS_AND_FLOW_MASK         0xF0000000
 112 #define V6_MCAST                        0xFF000000
 113 #define V6_LINKLOCAL                    0xFE800000

 115 #define IPV6_FLOW_TCLASS(x)             (((x) & IPV6_FLOWINFO_TCLASS) >> 20)
 116 #define IPV6_TCLASS_FLOW(f, c)          (((f) & ~IPV6_FLOWINFO_TCLASS) |\
 117                                         ((c) << 20))
 118 #else
 119 #define IPV6_DEFAULT_VERS_AND_FLOW      0x00000060
 120 #define IPV6_VERS_AND_FLOW_MASK         0x000000F0
 121 #define V6_MCAST                        0x000000FF
 122 #define V6_LINKLOCAL                    0x000080FE

 124 #define IPV6_FLOW_TCLASS(x)             ((((x) & 0xf000U) >> 12) |\
 125                                         (((x) & 0xf) << 4))
 126 #define IPV6_TCLASS_FLOW(f, c)          (((f) & ~IPV6_FLOWINFO_TCLASS) |\
 127                                         ((((c) & 0xf) << 12) |\
```

```
 128                                        (((c) & 0xf0) >> 4)))
 129 #endif

 131 /*
 132  * UTILITY MACROS FOR ADDRESSES.
 133  */

 135 /*
 136  * Convert an IPv4 address mask to an IPv6 mask.   Pad with 1-bits.
 137  */
 138 #define V4MASK_TO_V6(v4, v6)     ((v6).s6_addr32[0] = 0xffffffffUL,      \
 139                                   (v6).s6_addr32[1] = 0xffffffffUL,      \
 140                                   (v6).s6_addr32[2] = 0xffffffffUL,      \
 141                                   (v6).s6_addr32[3] = (v4))
 143 /*
 144  * Convert aligned IPv4-mapped IPv6 address into an IPv4 address.
 145  * Note: We use "v6" here in definition of macro instead of "(v6)"
 146  * Not possible to use "(v6)" here since macro is used with struct
 147  * field names as arguments.
 148  */
 149 #define V4_PART_OF_V6(v6)        v6.s6_addr32[3]

 151 #ifdef _BIG_ENDIAN
 152 #define V6_OR_V4_INADDR_ANY(a)  ((a).s6_addr32[3] == 0 &&              \
 153                                  ((a).s6_addr32[2] == 0xffffU ||       \
 154                                  (a).s6_addr32[2] == 0) &&             \
 155                                  (a).s6_addr32[1] == 0 &&              \
 156                                  (a).s6_addr32[0] == 0)

 158 #else
 159 #define V6_OR_V4_INADDR_ANY(a)  ((a).s6_addr32[3] == 0 &&              \
 160                                  ((a).s6_addr32[2] == 0xffff0000U ||   \
 161                                  (a).s6_addr32[2] == 0) &&             \
 162                                  (a).s6_addr32[1] == 0 &&              \
 163                                  (a).s6_addr32[0] == 0)
 164 #endif /* _BIG_ENDIAN */

 166 /* IPv4-mapped CLASSD addresses */
 167 #ifdef _BIG_ENDIAN
 168 #define IN6_IS_ADDR_V4MAPPED_CLASSD(addr) \
 169         (((addr)->_S6_un._S6_u32[2] == 0x0000ffff) && \
 170         (CLASSD((addr)->_S6_un._S6_u32[3])) && \
 171         ((addr)->_S6_un._S6_u32[1] == 0) && \
 172         ((addr)->_S6_un._S6_u32[0] == 0))
 173 #else   /* _BIG_ENDIAN */
 174 #define IN6_IS_ADDR_V4MAPPED_CLASSD(addr) \
 175         (((addr)->_S6_un._S6_u32[2] == 0xffff0000U) && \
 176         (CLASSD((addr)->_S6_un._S6_u32[3])) && \
 177         ((addr)->_S6_un._S6_u32[1] == 0) && \
 178         ((addr)->_S6_un._S6_u32[0] == 0))
 179 #endif /* _BIG_ENDIAN */

 181 /* Clear an IPv6 addr */
 182 #define V6_SET_ZERO(a)           ((a).s6_addr32[0] = 0,                \
 183                                   (a).s6_addr32[1] = 0,                \
 184                                   (a).s6_addr32[2] = 0,                \
 185                                   (a).s6_addr32[3] = 0)

 187 /* Mask comparison: is IPv6 addr a, and'ed with mask m, equal to addr b? */
 188 #define V6_MASK_EQ(a, m, b)                                            \
 189         ((((a).s6_addr32[0] & (m).s6_addr32[0]) == (b).s6_addr32[0]) && \
 190         (((a).s6_addr32[1] & (m).s6_addr32[1]) == (b).s6_addr32[1]) && \
 191         (((a).s6_addr32[2] & (m).s6_addr32[2]) == (b).s6_addr32[2]) && \
 192         (((a).s6_addr32[3] & (m).s6_addr32[3]) == (b).s6_addr32[3]))
```

```
 194 #define V6_MASK_EQ_2(a, m, b)                                          \
 195         ((((a).s6_addr32[0] & (m).s6_addr32[0]) ==                     \
 196             ((b).s6_addr32[0]  & (m).s6_addr32[0])) &&                 \
 197         (((a).s6_addr32[1] & (m).s6_addr32[1]) ==                      \
 198             ((b).s6_addr32[1]  & (m).s6_addr32[1])) &&                 \
 199         (((a).s6_addr32[2] & (m).s6_addr32[2]) ==                      \
 200             ((b).s6_addr32[2]  & (m).s6_addr32[2])) &&                 \
 201         (((a).s6_addr32[3] & (m).s6_addr32[3]) ==                      \
 202             ((b).s6_addr32[3]  & (m).s6_addr32[3])))

 204 /* Copy IPv6 address (s), logically and'ed with mask (m), into (d) */
 205 #define V6_MASK_COPY(s, m, d)                                          \
 206         ((d).s6_addr32[0] = (s).s6_addr32[0] & (m).s6_addr32[0],       \
 207         (d).s6_addr32[1] = (s).s6_addr32[1] & (m).s6_addr32[1],        \
 208         (d).s6_addr32[2] = (s).s6_addr32[2] & (m).s6_addr32[2],        \
 209         (d).s6_addr32[3] = (s).s6_addr32[3] & (m).s6_addr32[3])

 211 #define ILL_FRAG_HASH_V6(v6addr, i)                                    \
 212         ((ntohl((v6addr).s6_addr32[3]) ^ (i ^ (i >> 8))) %             \
 213                                         ILL_FRAG_HASH_TBL_COUNT)


 216 /*
 217  * GLOBAL EXTERNALS
 218  */
 219 extern const in6_addr_t ipv6_all_ones;
 220 extern const in6_addr_t ipv6_all_zeros;
 221 extern const in6_addr_t ipv6_loopback;
 222 extern const in6_addr_t ipv6_all_hosts_mcast;
 223 extern const in6_addr_t ipv6_all_rtrs_mcast;
 224 extern const in6_addr_t ipv6_all_v2rtrs_mcast;
 225 extern const in6_addr_t ipv6_solicited_node_mcast;
 226 extern const in6_addr_t ipv6_unspecified_group;

 228 /*
 229  * FUNCTION PROTOTYPES
 230  */

 232 extern void     convert2ascii(char *buf, const in6_addr_t *addr);
  25 extern char     *inet_ntop(int, const void *, char *, int);
  26 extern int      inet_pton(int, char *, void *);
 233 extern void     icmp_param_problem_nexthdr_v6(mblk_t *, boolean_t,
 234     ip_recv_attr_t *);
 235 extern void     icmp_pkt2big_v6(mblk_t *, uint32_t, boolean_t,
 236     ip_recv_attr_t *);
 237 extern void     icmp_time_exceeded_v6(mblk_t *, uint8_t, boolean_t,
 238     ip_recv_attr_t *);
 239 extern void     icmp_unreachable_v6(mblk_t *, uint8_t, boolean_t,
 240     ip_recv_attr_t *);
 241 extern mblk_t   *icmp_inbound_v6(mblk_t *, ip_recv_attr_t *);
 242 extern void     icmp_inbound_error_fanout_v6(mblk_t *, icmp6_t *,
 243     ip_recv_attr_t *);
 244 extern void     icmp_update_out_mib_v6(ill_t *, icmp6_t *);

 246 extern boolean_t conn_wantpacket_v6(conn_t *, ip_recv_attr_t *, ip6_t *);

 248 extern in6addr_scope_t  ip_addr_scope_v6(const in6_addr_t *);
 249 extern void     ip_build_hdrs_v6(uchar_t *, uint_t, const ip_pkt_t *, uint8_t,
 250     uint32_t);
 251 extern void     ip_fanout_udp_multi_v6(mblk_t *, ip6_t *, uint16_t, uint16_t,
 252     ip_recv_attr_t *);
 253 extern void     ip_fanout_send_icmp_v6(mblk_t *, uint_t, uint8_t,
 254     ip_recv_attr_t *);
 255 extern void     ip_fanout_proto_v6(mblk_t *, ip6_t *, ip_recv_attr_t *);
 256 extern int      ip_find_hdr_v6(mblk_t *, ip6_t *, boolean_t, ip_pkt_t *,
 257     uint8_t *);
```

```
258 extern in6_addr_t ip_get_dst_v6(ip6_t *, const mblk_t *, boolean_t *);
259 extern ip6_rthdr_t      *ip_find_rthdr_v6(ip6_t *, uint8_t *);
260 extern boolean_t        ip_hdr_length_nexthdr_v6(mblk_t *, ip6_t *,
261     uint16_t *, uint8_t **);
262 extern int       ip_hdr_length_v6(mblk_t *, ip6_t *);
263 extern uint32_t ip_massage_options_v6(ip6_t *, ip6_rthdr_t *, netstack_t *);
264 extern void      ip_forward_xmit_v6(nce_t *, mblk_t *, ip6_t *, ip_recv_attr_t *,
265     uint32_t, uint32_t);
266 extern mblk_t   *ip_fraghdr_add_v6(mblk_t *, uint32_t, ip_xmit_attr_t *);
267 extern int       ip_fragment_v6(mblk_t *, nce_t *, iaflags_t, uint_t, uint32_t,
268     uint32_t, zoneid_t, zoneid_t, pfirepostfrag_t postfragfn,
269     uintptr_t *ixa_cookie);
270 extern int       ip_process_options_v6(mblk_t *, ip6_t *,
271     uint8_t *, uint_t, uint8_t, ip_recv_attr_t *);
272 extern void      ip_process_rthdr(mblk_t *, ip6_t *, ip6_rthdr_t *,
273     ip_recv_attr_t *);
274 extern int       ip_total_hdrs_len_v6(const ip_pkt_t *);
275 extern mblk_t   *ipsec_early_ah_v6(mblk_t *, ip_recv_attr_t *);
276 extern int       ipsec_ah_get_hdr_size_v6(mblk_t *, boolean_t);
277 extern void      ip_send_potential_redirect_v6(mblk_t *, ip6_t *, ire_t *,
278     ip_recv_attr_t *);
279 extern void      ip_rput_v6(queue_t *, mblk_t *);
280 extern mblk_t   *mld_input(mblk_t *, ip_recv_attr_t *);
281 extern void      mld_joingroup(ilm_t *);
282 extern void      mld_leavegroup(ilm_t *);
283 extern void      mld_timeout_handler(void *);

285 extern void      pr_addr_dbg(char *, int, const void *);
286 extern void      *ip6_kstat_init(netstackid_t, ip6_stat_t *);
287 extern void      ip6_kstat_fini(netstackid_t, kstat_t *);
288 extern size_t   ip6_get_src_preferences(ip_xmit_attr_t *, uint32_t *);
289 extern int       ip6_set_src_preferences(ip_xmit_attr_t *, uint32_t);

291 #endif  /* _KERNEL */

293 #ifdef  __cplusplus
294 }
_____unchanged_portion_omitted_
```

```
*********************************************************
   42357 Wed Sep 26 12:51:48 2012
new/usr/src/uts/common/io/scsi/adapters/iscsi/iscsi.h
inet_pton
*********************************************************
   1 /*
   2  * CDDL HEADER START
   3  *
   4  * The contents of this file are subject to the terms of the
   5  * Common Development and Distribution License (the "License").
   6  * You may not use this file except in compliance with the License.
   7  *
   8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
   9  * or http://www.opensolaris.org/os/licensing.
  10  * See the License for the specific language governing permissions
  11  * and limitations under the License.
  12  *
  13  * When distributing Covered Code, include this CDDL HEADER in each
  14  * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
  15  * If applicable, add the following below this CDDL HEADER, with the
  16  * fields enclosed by brackets "[]" replaced with your own identifying
  17  * information: Portions Copyright [yyyy] [name of copyright owner]
  18  *
  19  * CDDL HEADER END
  20  */
  21 /*
  22  * Copyright 2000 by Cisco Systems, Inc.  All rights reserved.
  23  * Copyright (c) 2008, 2010, Oracle and/or its affiliates. All rights reserved.
  24  * Copyright 2012 Nexenta Systems, Inc. All rights reserved.
  25 #endif /* ! codereview */
  26  */

  28 #ifndef _ISCSI_H
  29 #define _ISCSI_H

  31 /*
  32  * Block comment which describes the contents of this file.
  33  */

  35 #ifdef __cplusplus
  36 extern "C" {
  37 #endif

  39 #include <sys/scsi/scsi.h>
  40 #include <sys/ddi.h>
  41 #include <sys/sunddi.h>
  42 #include <sys/socket.h>
  43 #include <sys/kstat.h>
  44 #include <sys/sunddi.h>
  45 #include <sys/sunmdi.h>
  46 #include <sys/mdi_impldefs.h>
  47 #include <sys/time.h>
  48 #include <sys/nvpair.h>
  49 #include <sys/sdt.h>

  51 #include <sys/iscsi_protocol.h>
  52 #include <sys/scsi/adapters/iscsi_if.h>
  53 #include <iscsiAuthClient.h>
  54 #include <iscsi_stats.h>
  55 #include <iscsi_thread.h>
  56 #include <sys/idm/idm.h>
  57 #include <sys/idm/idm_conn_sm.h>
  58 #include <nvfile.h>
  59 #include <inet/ip.h>
  60 #endif /* ! codereview */
```

```
  62 #ifndef MIN
  63 #define MIN(a, b) ((a) < (b) ? (a) : (b))
  64 #endif

  66 #ifndef TRUE
  67 #define TRUE 1
  68 #endif

  70 #ifndef FALSE
  71 #define FALSE 0
  72 #endif

  74 #define LOGIN_PDU_BUFFER_SIZE   (16 * 1024)     /* move somewhere else */

  76 extern boolean_t iscsi_conn_logging;
  77 extern boolean_t iscsi_io_logging;
  78 extern boolean_t iscsi_login_logging;
  79 extern boolean_t iscsi_logging;
  80 extern boolean_t iscsi_sess_logging;
  81 #define ISCSI_CONN_LOG  if (iscsi_conn_logging) cmn_err
  82 #define ISCSI_IO_LOG    if (iscsi_io_logging) cmn_err
  83 #define ISCSI_LOGIN_LOG if (iscsi_login_logging) cmn_err
  84 #define ISCSI_LOG       if (iscsi_logging) cmn_err
  85 #define ISCSI_SESS_LOG  if (iscsi_sess_logging) cmn_err

  87 /*
  88  * Name Format of the different Task Queues
  89  */
  90 #define ISCSI_SESS_IOTH_NAME_FORMAT             "io_thrd_%d.%d"
  91 #define ISCSI_SESS_WD_NAME_FORMAT               "wd_thrd_%d.%d"
  92 #define ISCSI_SESS_LOGIN_TASKQ_NAME_FORMAT      "login_taskq_%d.%d"
  93 #define ISCSI_SESS_ENUM_TASKQ_NAME_FORMAT       "enum_taskq_%d.%d"
  94 #define ISCSI_CONN_CN_TASKQ_NAME_FORMAT         "conn_cn_taskq_%d.%d.%d"
  95 #define ISCSI_CONN_RXTH_NAME_FORMAT             "rx_thrd_%d.%d.%d"
  96 #define ISCSI_CONN_TXTH_NAME_FORMAT             "tx_thrd_%d.%d.%d"

  98 /*
  99  * The iSCSI driver will not build scatter/gather lists (iovec) longer
 100  * than the value defined here. Asserts have been include in the code
 101  * to check.
 102  */
 103 #define ISCSI_MAX_IOVEC         5

 105 #define ISCSI_DEFAULT_MAX_STORM_DELAY           32

 107 /*
 108  * The SNDBUF and RCVBUF size parameters for the sockets are just a
 109  * guess for the time being (I think it is the values used by CISCO
 110  * or UNH).  Testing will have to be done to figure * out the impact
 111  * of these values on performance.
 112  */
 113 #define ISCSI_SOCKET_SNDBUF_SIZE                (256 * 1024)
 114 #define ISCSI_SOCKET_RCVBUF_SIZE                (256 * 1024)
 115 #define ISCSI_TCP_NODELAY_DEFAULT               0
 116 #define ISCSI_TCP_CNOTIFY_THRESHOLD_DEFAULT     2000
 117 #define ISCSI_TCP_CABORT_THRESHOLD_DEFAULT      10000
 118 #define ISCSI_TCP_ABORT_THRESHOLD_DEFAULT       (30 * 1000) /* milliseconds */
 119 #define ISNS_TCP_ABORT_THRESHOLD_DEFAULT        (3 * 1000) /* milliseconds */

 121 /* Default values for tunable parameters */
 122 #define ISCSI_DEFAULT_RX_TIMEOUT_VALUE          60
 123 #define ISCSI_DEFAULT_CONN_DEFAULT_LOGIN_MAX    180
 124 #define ISCSI_DEFAULT_LOGIN_POLLING_DELAY       60

 126 /*
 127  * Convenient short hand defines
```

```
 128  */
 129 #define TARGET_PROP      "target"
 130 #define LUN_PROP         "lun"
 131 #define MDI_GUID         "wwn"
 132 #define NDI_GUID         "client-guid"

 134 #define ISCSI_SIG_CMD    0x11111111
 135 #define ISCSI_SIG_LUN    0x22222222
 136 #define ISCSI_SIG_CONN   0x33333333
 137 #define ISCSI_SIG_SESS   0x44444444
 138 #define ISCSI_SIG_HBA    0x55555555

 140 #define SENDTARGETS_DISCOVERY   "SENDTARGETS_DISCOVERY"

 142 #define ISCSI_LUN_MASK_MSB      0x00003f00
 143 #define ISCSI_LUN_MASK_LSB      0x000000ff
 144 #define ISCSI_LUN_MASK          (ISCSI_LUN_MASK_MSB | ISCSI_LUN_MASK_LSB)
 145 #define ISCSI_LUN_BYTE_COPY(lun, report_lun_data) \
 146         lun[0] = (report_lun_data & ISCSI_LUN_MASK_MSB) >> 8; \
 147         lun[1] = (report_lun_data & ISCSI_LUN_MASK_LSB);
 148 /*
 149  * Not defined by iSCSI, but used in the login code to
 150  * determine when to send the initial Login PDU
 151  */
 152 #define ISCSI_INITIAL_LOGIN_STAGE       -1

 154 typedef enum iscsi_status {
 155         /* Success */
 156         ISCSI_STATUS_SUCCESS = 0,
 157         /* Driver / Kernel / Code error */
 158         ISCSI_STATUS_INTERNAL_ERROR,
 159         /* ITT table is already full, unable to reserve slot */
 160         ISCSI_STATUS_ITT_TABLE_FULL,
 161         /* Login on connection failed */
 162         ISCSI_STATUS_LOGIN_FAILED,
 163         /* No connections are in the LOGGED_IN state */
 164         ISCSI_STATUS_NO_CONN_LOGGED_IN,
 165         /* TCP Transfer Error */
 166         ISCSI_STATUS_TCP_TX_ERROR,
 167         /* TCP Receive Error */
 168         ISCSI_STATUS_TCP_RX_ERROR,
 169         /* iSCSI packet RCV timeout */
 170         ISCSI_STATUS_RX_TIMEOUT,
 171         /* iSCSI Header Digest CRC error */
 172         ISCSI_STATUS_HEADER_DIGEST_ERROR,
 173         /* iSCSI Data Digest CRC error */
 174         ISCSI_STATUS_DATA_DIGEST_ERROR,
 175         /* kmem_alloc failure */
 176         ISCSI_STATUS_ALLOC_FAILURE,
 177         /* cmd (tran_abort/reset) failed */
 178         ISCSI_STATUS_CMD_FAILED,
 179         /* iSCSI protocol error */
 180         ISCSI_STATUS_PROTOCOL_ERROR,
 181         /* iSCSI protocol version mismatch */
 182         ISCSI_STATUS_VERSION_MISMATCH,
 183         /* iSCSI login negotiation failed */
 184         ISCSI_STATUS_NEGO_FAIL,
 185         /* iSCSI login authentication failed */
 186         ISCSI_STATUS_AUTHENTICATION_FAILED,
 187         /* iSCSI login redirection failed */
 188         ISCSI_STATUS_REDIRECTION_FAILED,
 189         /* iSCSI uscsi status failure */
 190         ISCSI_STATUS_USCSI_FAILED,
 191         /* data received would have overflowed given buffer */
 192         ISCSI_STATUS_DATA_OVERFLOW,
 193         /* session/connection needs to shutdown */
```

```
 194         ISCSI_STATUS_SHUTDOWN,
 195         /* logical unit in use */
 196         ISCSI_STATUS_BUSY,
 197         /* Login on connection failed, retries exceeded */
 198         ISCSI_STATUS_LOGIN_TIMED_OUT,
 199         /* iSCSI login tpgt negotiation failed */
 200         ISCSI_STATUS_LOGIN_TPGT_NEGO_FAIL
 201 } iscsi_status_t;
 202 #define ISCSI_SUCCESS(status) (status == ISCSI_STATUS_SUCCESS)

 204 /* SNA32 check value used on increment of CmdSn values */
 205 #define ISCSI_SNA32_CHECK 2147483648UL /* 2**31 */

 207 /*
 208  * This is the maximum number of commands that can be outstanding
 209  * on a iSCSI session at anyone point in time.
 210  */
 211 #define ISCSI_CMD_TABLE_SIZE            1024

 213 /* Used on connections thread create of receiver thread */
 214 extern pri_t minclsyspri;

 216 /*
 217  * Callers of iscsid_config_one/all must hold this
 218  * semaphore across the calls.  Otherwise a ndi_devi_enter()
 219  * deadlock in the DDI layer may occur.
 220  */
 221 extern ksema_t iscsid_config_semaphore;

 223 extern kmutex_t iscsi_oid_mutex;
 224 extern uint32_t iscsi_oid;
 225 extern void *iscsi_state;

 227 /*
 228  * NOP delay is used to send a iSCSI NOP (ie. ping) across the
 229  * wire to see if the target is still alive.  NOPs are only
 230  * sent when the RX thread hasn't received anything for the
 231  * below amount of time.
 232  */
 233 #define ISCSI_DEFAULT_NOP_DELAY                 5 /* seconds */
 234 extern int      iscsi_nop_delay;
 235 /*
 236  * If we haven't received anything in a specified period of time
 237  * we will stop accepting IO via tran start.  This will enable
 238  * upper level drivers to see we might be having a problem and
 239  * in the case of scsi_vhci will start to route IO down a better
 240  * path.
 241  */
 242 #define ISCSI_DEFAULT_RX_WINDOW                 20 /* seconds */
 243 extern int      iscsi_rx_window;
 244 /*
 245  * If we haven't received anything in a specified period of time
 246  * we will stop accepting IO via tran start.  This the max limit
 247  * when encountered we will start returning a fatal error.
 248  */
 249 #define ISCSI_DEFAULT_RX_MAX_WINDOW             180 /* seconds */
 250 extern int      iscsi_rx_max_window;

 252 /*
 253  * During iscsi boot, if the boot session has been created, the
 254  * initiator hasn't changed the boot lun to be online, we will wait
 255  * 180s here for lun online by default.
 256  */
 257 #define ISCSI_BOOT_DEFAULT_MAX_DELAY            180 /* seconds */
 258 /*
 259  * +-------------------------------------------------------------------+
```

```
 260  * | iSCSI Driver Structures                                              |
 261  * +----------------------------------------------------------------------+
 262  */

 264 /*
 265  * iSCSI Auth Information
 266  */
 267 typedef struct iscsi_auth {
 268         IscsiAuthStringBlock    auth_recv_string_block;
 269         IscsiAuthStringBlock    auth_send_string_block;
 270         IscsiAuthLargeBinary    auth_recv_binary_block;
 271         IscsiAuthLargeBinary    auth_send_binary_block;
 272         IscsiAuthClient         auth_client_block;
 273         int                     num_auth_buffers;
 274         IscsiAuthBufferDesc     auth_buffers[5];

 276         /*
 277          * To indicate if bi-directional authentication is enabled.
 278          * 0 means uni-directional authentication.
 279          * 1 means bi-directional authentication.
 280          */
 281         int                     bidirectional_auth;

 283         /* Initiator's authentication information. */
 284         char                    username[iscsiAuthStringMaxLength];
 285         uint8_t                 password[iscsiAuthStringMaxLength];
 286         int                     password_length;

 288         /* Target's authentication information. */
 289         char                    username_in[iscsiAuthStringMaxLength];
 290         uint8_t                 password_in[iscsiAuthStringMaxLength];
 291         int                     password_length_in;
 292 } iscsi_auth_t;

 294 /*
 295  * iSCSI Task
 296  */
 297 typedef struct iscsi_task {
 298         void                    *t_arg;
 299         boolean_t               t_blocking;
 300         uint32_t                t_event_count;
 301 } iscsi_task_t;

 303 /*
 304  * These are all the iscsi_cmd types that we use to track our
 305  * commands between queues and actions.
 306  */
 307 typedef enum iscsi_cmd_type {
 308         ISCSI_CMD_TYPE_SCSI = 1,        /* scsi cmd */
 309         ISCSI_CMD_TYPE_NOP,             /* nop / ping */
 310         ISCSI_CMD_TYPE_ABORT,           /* abort */
 311         ISCSI_CMD_TYPE_RESET,           /* reset */
 312         ISCSI_CMD_TYPE_LOGOUT,          /* logout */
 313         ISCSI_CMD_TYPE_LOGIN,           /* login */
 314         ISCSI_CMD_TYPE_TEXT             /* text */
 315 } iscsi_cmd_type_t;

 317 /*
 318  * iscsi_cmd_state - (reference iscsi_cmd.c for state diagram)
 319  */
 320 typedef enum iscsi_cmd_state {
 321         ISCSI_CMD_STATE_FREE = 0,
 322         ISCSI_CMD_STATE_PENDING,
 323         ISCSI_CMD_STATE_ACTIVE,
 324         ISCSI_CMD_STATE_ABORTING,
 325         ISCSI_CMD_STATE_IDM_ABORTING,
```

```
 326         ISCSI_CMD_STATE_COMPLETED,
 327         ISCSI_CMD_STATE_MAX
 328 } iscsi_cmd_state_t;

 330 #ifdef ISCSI_CMD_SM_STRINGS
 331 static const char *iscsi_cmd_state_names[ISCSI_CMD_STATE_MAX+1] = {
 332         "ISCSI_CMD_STATE_FREE",
 333         "ISCSI_CMD_STATE_PENDING",
 334         "ISCSI_CMD_STATE_ACTIVE",
 335         "ISCSI_CMD_STATE_ABORTING",
 336         "ISCSI_CMD_STATE_IDM_ABORTING",
 337         "ISCSI_CMD_STATE_COMPLETED",
 338         "ISCSI_CMD_STATE_MAX"
 339 };
 340 #endif

 342 /*
 343  * iscsi command events
 344  */
 345 typedef enum iscsi_cmd_event {
 346         ISCSI_CMD_EVENT_E1 = 0,
 347         ISCSI_CMD_EVENT_E2,
 348         ISCSI_CMD_EVENT_E3,
 349         ISCSI_CMD_EVENT_E4,
 350         ISCSI_CMD_EVENT_E6,
 351         ISCSI_CMD_EVENT_E7,
 352         ISCSI_CMD_EVENT_E8,
 353         ISCSI_CMD_EVENT_E9,
 354         ISCSI_CMD_EVENT_E10,
 355         ISCSI_CMD_EVENT_MAX
 356 } iscsi_cmd_event_t;

 358 #ifdef ISCSI_CMD_SM_STRINGS
 359 static const char *iscsi_cmd_event_names[ISCSI_CMD_EVENT_MAX+1] = {
 360         "ISCSI_CMD_EVENT_E1",
 361         "ISCSI_CMD_EVENT_E2",
 362         "ISCSI_CMD_EVENT_E3",
 363         "ISCSI_CMD_EVENT_E4",
 364         "ISCSI_CMD_EVENT_E6",
 365         "ISCSI_CMD_EVENT_E7",
 366         "ISCSI_CMD_EVENT_E8",
 367         "ISCSI_CMD_EVENT_E9",
 368         "ISCSI_CMD_EVENT_E10",
 369         "ISCSI_CMD_EVENT_MAX"
 370 };
 371 #endif

 373 /*
 374  * iscsi text command stages - these stages are used by iSCSI text
 375  * processing to manage long resonses.
 376  */
 377 typedef enum iscsi_cmd_text_stage {
 378         ISCSI_CMD_TEXT_INITIAL_REQ = 0,
 379         ISCSI_CMD_TEXT_CONTINUATION,
 380         ISCSI_CMD_TEXT_FINAL_RSP
 381 } iscsi_cmd_text_stage_t;

 383 /*
 384  * iscsi cmd misc flags - bitwise applicable
 385  */
 386 #define ISCSI_CMD_MISCFLAG_INTERNAL     0x1
 387 #define ISCSI_CMD_MISCFLAG_FREE         0x2
 388 #define ISCSI_CMD_MISCFLAG_STUCK        0x4
 389 #define ISCSI_CMD_MISCFLAG_XARQ         0x8
 390 #define ISCSI_CMD_MISCFLAG_SENT         0x10
 391 #define ISCSI_CMD_MISCFLAG_FLUSH        0x20
```

```
 393 /*
 394  * 1/2 of a 32 bit number, used for checking CmdSN
 395  * wrapped.
 396  */
 397 #define ISCSI_CMD_SN_WRAP              0x80000000

 399 #define ISCSI_CMD_PKT_STAT_INIT        0

 401 /*
 402  * iSCSI cmd/pkt Structure
 403  */
 404 typedef struct iscsi_cmd {
 405         uint32_t                cmd_sig;
 406         struct iscsi_cmd        *cmd_prev;
 407         struct iscsi_cmd        *cmd_next;
 408         struct iscsi_conn       *cmd_conn;

 410         iscsi_cmd_type_t        cmd_type;
 411         iscsi_cmd_state_t       cmd_state;
 412         iscsi_cmd_state_t       cmd_prev_state;
 413         clock_t                 cmd_lbolt_pending;
 414         clock_t                 cmd_lbolt_active;
 415         clock_t                 cmd_lbolt_aborting;
 416         clock_t                 cmd_lbolt_idm_aborting;
 417         clock_t                 cmd_lbolt_timeout;
 418         uint8_t                 cmd_misc_flags;
 419         idm_task_t              *cmd_itp;

 421         union {
 422                         /* ISCSI_CMD_TYPE_SCSI */
 423                 struct {
 424                         idm_buf_t               *ibp_ibuf;
 425                         idm_buf_t               *ibp_obuf;
 426                         struct scsi_pkt         *pkt;
 427                         struct buf              *bp;
 428                         int                     cmdlen;
 429                         int                     statuslen;
 430                         size_t                  data_transferred;

 432                         uint32_t                lun;

 434                         /*
 435                          * If SCSI_CMD_TYPE is in ABORTING_STATE
 436                          * then the abort_icmdp field will be a pointer
 437                          * to the abort command chasing this one.
 438                          */
 439                         struct iscsi_cmd        *abort_icmdp;
 440                         /*
 441                          * pointer to the r2t associated with this
 442                          * command (if any)
 443                          */
 444                         struct iscsi_cmd        *r2t_icmdp;
 445                         /*
 446                          * It will be true if this command has
 447                          * another R2T to handle.
 448                          */
 449                         boolean_t               r2t_more;
 450                         /*
 451                          * It is used to record pkt_statistics temporarily.
 452                          */
 453                         uint_t                  pkt_stat;
 454                 } scsi;
 455                         /* ISCSI_CMD_TYPE_ABORT */
 456                 struct {
 457                         /* pointer to original iscsi_cmd, for abort */
```

```
 458                         struct iscsi_cmd        *icmdp;
 459                 } abort;
 460                         /* ISCSI_CMD_TYPE_RESET */
 461                 struct {
 462                         int                     level;
 463                         uint8_t                 response;
 464                 } reset;
 465                         /* ISCSI_CMD_TYPE_NOP */
 466                 struct {
 467                         int rsvd;
 468                 } nop;
 469                         /* ISCSI_CMD_TYPE_R2T */
 470                 struct {
 471                         struct iscsi_cmd        *icmdp;
 472                         uint32_t                offset;
 473                         uint32_t                length;
 474                 } r2t;
 475                         /* ISCSI_CMD_TYPE_LOGIN */
 476                 struct {
 477                         int rvsd;
 478                 } login;
 479                         /* ISCSI_CMD_TYPE_LOGOUT */
 480                 struct {
 481                         int rsvd;
 482                 } logout;
 483                         /* ISCSI_CMD_TYPE_TEXT */
 484                 struct {
 485                         char                    *buf;
 486                         int                     buf_len;
 487                         uint32_t                offset;
 488                         uint32_t                data_len;
 489                         uint32_t                total_rx_len;
 490                         uint32_t                ttt;
 491                         uint8_t                 lun[8];
 492                         iscsi_cmd_text_stage_t  stage;
 493                 } text;
 494         } cmd_un;

 496         struct iscsi_lun        *cmd_lun; /* associated lun */

 498         uint32_t                cmd_itt;
 499         uint32_t                cmd_ttt;

 501         /*
 502          * If a data digest error is seem on a data pdu.  This flag
 503          * will get set.  We don't abort the cmd immediately because
 504          * we want to read in all the data to get it out of the
 505          * stream.  Once the completion for the cmd is received we
 506          * we will abort the cmd and state no sense data was available.
 507          */
 508         boolean_t               cmd_crc_error_seen;

 510         /*
 511          * Used to block and wake up caller until action is completed.
 512          * This is for ABORT, RESET, and PASSTHRU cmds.
 513          */
 514         int                     cmd_result;
 515         int                     cmd_completed;
 516         kmutex_t                cmd_mutex;
 517         kcondvar_t              cmd_completion;

 519         idm_pdu_t               cmd_pdu;

 521         sm_audit_buf_t          cmd_state_audit;

 523         uint32_t                cmd_sn;
```

```
 524 } iscsi_cmd_t;


 527 /*
 528  * iSCSI LUN Structure
 529  */
 530 typedef struct iscsi_lun {
 531         uint32_t                lun_sig;
 532         int                     lun_state;

 534         struct iscsi_lun        *lun_next;      /* next lun on this sess. */
 535         struct iscsi_sess       *lun_sess;      /* parent sess. for lun */
 536         dev_info_t              *lun_dip;
 537         mdi_pathinfo_t          *lun_pip;

 539         uint16_t                lun_num;        /* LUN */
 540         uint8_t                 lun_addr_type;  /* LUN addressing type */
 541         uint32_t                lun_oid;        /* OID */
 542         char                    *lun_guid;      /* GUID */
 543         int                     lun_guid_size;  /* GUID allocation size */
 544         char                    *lun_addr;      /* sess,lun */
 545         time_t                  lun_time_online;

 547         uchar_t                 lun_cap;        /* bitmap of scsi caps */

 549         uchar_t                 lun_vid[ISCSI_INQ_VID_BUF_LEN]; /* Vendor ID */
 550         uchar_t                 lun_pid[ISCSI_INQ_PID_BUF_LEN]; /* Product ID */

 552         uchar_t                 lun_type;
 553 } iscsi_lun_t;

 555 #define ISCSI_LUN_STATE_CLEAR       0               /* used to clear all states */
 556 #define ISCSI_LUN_STATE_OFFLINE     1
 557 #define ISCSI_LUN_STATE_ONLINE      2
 558 #define ISCSI_LUN_STATE_INVALID     4               /* offline failed */
 559 #define ISCSI_LUN_STATE_BUSY        8               /* logic unit is in reset */

 561 #define ISCSI_LUN_CAP_RESET         0x01

 563 #define ISCSI_SCSI_RESET_SENSE_CODE 0x29
 564 #define ISCSI_SCSI_LUNCHANGED_CODE      0x3f

 566 #define ISCSI_SCSI_LUNCHANGED_ASCQ      0x0e

 568 /*
 569  *
 570  *
 571  */
 572 typedef struct iscsi_queue {
 573         iscsi_cmd_t     *head;
 574         iscsi_cmd_t     *tail;
 575         int             count;
 576         kmutex_t        mutex;
 577 } iscsi_queue_t;

 579 #define ISCSI_CONN_DEFAULT_LOGIN_MIN            0
 580 #define ISCSI_CONN_DEFAULT_LOGIN_REDIRECT       10

 582 /* iSCSI tunable Parameters */
 583 typedef struct iscsi_tunable_params {
 584         int             recv_login_rsp_timeout; /* range: 0 - 60*60 */
 585         int             conn_login_max;         /* range: 0 - 60*60 */
 586         int             polling_login_delay;    /* range: 0 - 60*60 */
 587 } iscsi_tunable_params_t;

 589 typedef union iscsi_sockaddr {
```

```
 590         struct sockaddr         sin;
 591         struct sockaddr_in      sin4;
 592         struct sockaddr_in6     sin6;
 593 } iscsi_sockaddr_t;

 595 #define SIZEOF_SOCKADDR(so)     ((so)->sa_family == AF_INET ? \
 596         sizeof (struct sockaddr_in) : sizeof (struct sockaddr_in6))

 598 typedef enum {
 599         LOGIN_START,
 600         LOGIN_READY,
 601         LOGIN_TX,
 602         LOGIN_RX,
 603         LOGIN_ERROR,
 604         LOGIN_DONE,
 605         LOGIN_FFP,
 606         LOGIN_MAX
 607 } iscsi_login_state_t;

 609 #ifdef ISCSI_LOGIN_STATE_NAMES
 610 static const char *iscsi_login_state_names[LOGIN_MAX+1] = {
 611         "LOGIN_START",
 612         "LOGIN_READY",
 613         "LOGIN_TX",
 614         "LOGIN_RX",
 615         "LOGIN_ERROR",
 616         "LOGIN_DONE",
 617         "LOGIN_FFP",
 618         "LOGIN_MAX"
 619 };
 620 #endif

 622 /*
 623  * iscsi_conn_state
 624  */
 625 typedef enum iscsi_conn_state {
 626         ISCSI_CONN_STATE_UNDEFINED = 0,
 627         ISCSI_CONN_STATE_FREE,
 628         ISCSI_CONN_STATE_IN_LOGIN,
 629         ISCSI_CONN_STATE_LOGGED_IN,
 630         ISCSI_CONN_STATE_IN_LOGOUT,
 631         ISCSI_CONN_STATE_FAILED,
 632         ISCSI_CONN_STATE_POLLING,
 633         ISCSI_CONN_STATE_MAX
 634 } iscsi_conn_state_t;

 636 #ifdef ISCSI_ICS_NAMES
 637 static const char *iscsi_ics_name[ISCSI_CONN_STATE_MAX+1] = {
 638         "ISCSI_CONN_STATE_UNDEFINED",
 639         "ISCSI_CONN_STATE_FREE",
 640         "ISCSI_CONN_STATE_IN_LOGIN",
 641         "ISCSI_CONN_STATE_LOGGED_IN",
 642         "ISCSI_CONN_STATE_IN_LOGOUT",
 643         "ISCSI_CONN_STATE_FAILED",
 644         "ISCSI_CONN_STATE_POLLING",
 645         "ISCSI_CONN_STATE_MAX"
 646 };
 647 #endif

 649 #define ISCSI_CONN_STATE_FULL_FEATURE(state) \
 650         ((state == ISCSI_CONN_STATE_LOGGED_IN) || \
 651         (state == ISCSI_CONN_STATE_IN_LOGOUT))

 653 /*
 654  * iSCSI Connection Structure
 655  */
```

```
 656 typedef struct iscsi_conn {
 657         uint32_t                conn_sig;
 658         struct iscsi_conn       *conn_next;     /* next conn on this sess. */
 659         struct iscsi_sess       *conn_sess;     /* parent sess. for conn. */

 661         iscsi_conn_state_t      conn_state;     /* cur. conn. driver state */
 662         iscsi_conn_state_t      conn_prev_state; /* prev. conn. driver state */
 663         /* protects the session state and synchronizes the state machine */
 664         kmutex_t                conn_state_mutex;
 665         kcondvar_t              conn_state_change;
 666         boolean_t               conn_state_destroy;
 667         boolean_t               conn_state_ffp;
 668         boolean_t               conn_state_idm_connected;
 669         boolean_t               conn_async_logout;
 670         ddi_taskq_t             *conn_cn_taskq;

 672         idm_conn_t              *conn_ic;

 674         /* base connection information, may have been redirected */
 675         iscsi_sockaddr_t        conn_base_addr;

 677         /* current connection information, may have been redirected */
 678         iscsi_sockaddr_t        conn_curr_addr;

 680         boolean_t               conn_bound;
 681         iscsi_sockaddr_t        conn_bound_addr;

 683         uint32_t                conn_cid;       /* CID */
 684         uint32_t                conn_oid;       /* OID */

 686         int                     conn_current_stage;    /* iSCSI login stage */
 687         int                     conn_next_stage;       /* iSCSI login stage */
 688         int                     conn_partial_response;

 690         /*
 691          * The active queue contains iscsi_cmds that have already
 692          * been sent on this connection.  Any future responses to
 693          * these cmds require alligence to this connection.  If there
 694          * are issues with these cmds the command may need aborted
 695          * depending on the command state, and must be put back into
 696          * the session's pending queue or aborted.
 697          */
 698         iscsi_queue_t           conn_queue_active;
 699         iscsi_queue_t           conn_queue_idm_aborting;

 701         /* lbolt from the last receive, used for nop processing */
 702         clock_t                 conn_rx_lbolt;
 703         clock_t                 conn_nop_lbolt;

 705         iscsi_thread_t          *conn_tx_thread;

 707         /*
 708          * The expstatsn is the command status sn that is expected
 709          * next from the target.  Command status is carried on a number
 710          * of iSCSI PDUs (ex.  SCSI Cmd Response, SCSI Data IN with
 711          * S-Bit set, ...), not all PDUs.  If our expstatsn is different
 712          * than the received statsn.  Something got out of sync we need to
 713          * recover.
 714          */
 715         uint32_t                conn_expstatsn;
 716         uint32_t                conn_laststatsn;

 718         /* active login parameters */
 719         iscsi_login_params_t    conn_params;

 721         /* Statistics */
```

```
 722         struct {
 723                 kstat_t                 *ks;
 724                 iscsi_conn_stats_t      ks_data;
 725         } stats;

 727         /*
 728          * These fields are used to coordinate the asynchronous IDM
 729          * PDU operations with the synchronous login code.
 730          */
 731         kmutex_t                conn_login_mutex;
 732         kcondvar_t              conn_login_cv;
 733         iscsi_login_state_t     conn_login_state;
 734         iscsi_status_t          conn_login_status;
 735         iscsi_hdr_t             conn_login_resp_hdr;
 736         char                    *conn_login_data;
 737         int                     conn_login_datalen;
 738         int                     conn_login_max_data_length;

 740         /*
 741          * login min and max identify the amount of time
 742          * in lbolt that iscsi_start_login() should attempt
 743          * to log into a target portal.  The login will
 744          * delay until the min lbolt has been reached and
 745          * will end once max time has been reached.  These
 746          * values are normally set to the default but can
 747          * are also altered by async commands received from
 748          * the targetlogin.
 749          */
 750         clock_t                 conn_login_min;
 751         clock_t                 conn_login_max;
 752         sm_audit_buf_t          conn_state_audit;

 754         /* active tunable parameters */
 755         iscsi_tunable_params_t  conn_tunable_params;
 756         boolean_t               conn_timeout;
 757 } iscsi_conn_t;


 760 /*
 761  * iscsi_sess_state - (reference iscsi_sess.c for state diagram)
 762  */
 763 typedef enum iscsi_sess_state {
 764         ISCSI_SESS_STATE_FREE = 0,
 765         ISCSI_SESS_STATE_LOGGED_IN,
 766         ISCSI_SESS_STATE_FAILED,
 767         ISCSI_SESS_STATE_IN_FLUSH,
 768         ISCSI_SESS_STATE_FLUSHED,
 769         ISCSI_SESS_STATE_MAX
 770 } iscsi_sess_state_t;

 772 #ifdef ISCSI_SESS_SM_STRINGS
 773 static const char *iscsi_sess_state_names[ISCSI_SESS_STATE_MAX+1] = {
 774         "ISCSI_SESS_STATE_FREE",
 775         "ISCSI_SESS_STATE_LOGGED_IN",
 776         "ISCSI_SESS_STATE_FAILED",
 777         "ISCSI_SESS_STATE_IN_FLUSH",
 778         "ISCSI_SESS_STATE_FLUSHED",
 779         "ISCSI_SESS_STATE_MAX"
 780 };
 781 #endif

 783 #define ISCSI_SESS_STATE_FULL_FEATURE(state) \
 784         ((state == ISCSI_SESS_STATE_LOGGED_IN) || \
 785         (state == ISCSI_SESS_STATE_IN_FLUSH))
```

```
788 typedef enum iscsi_sess_event {
789         ISCSI_SESS_EVENT_N1 = 0,
790         ISCSI_SESS_EVENT_N3,
791         ISCSI_SESS_EVENT_N5,
792         ISCSI_SESS_EVENT_N6,
793         ISCSI_SESS_EVENT_N7,
794         ISCSI_SESS_EVENT_MAX
795 } iscsi_sess_event_t;

797 #ifdef ISCSI_SESS_SM_STRINGS
798 static const char *iscsi_sess_event_names[ISCSI_SESS_EVENT_MAX+1] = {
799         "ISCSI_SESS_EVENT_N1",
800         "ISCSI_SESS_EVENT_N3",
801         "ISCSI_SESS_EVENT_N5",
802         "ISCSI_SESS_EVENT_N6",
803         "ISCSI_SESS_EVENT_N7",
804         "ISCSI_SESS_EVENT_MAX"
805 };
806 #endif

808 typedef enum iscsi_sess_type {
809         ISCSI_SESS_TYPE_NORMAL = 0,
810         ISCSI_SESS_TYPE_DISCOVERY
811 } iscsi_sess_type_t;

813 #define SESS_ABORT_TASK_MAX_THREADS     1

815 /* Sun's initiator session ID */
816 #define ISCSI_SUN_ISID_0    0x40    /* ISID - EN format */
817 #define ISCSI_SUN_ISID_1    0x00    /* Sec B */
818 #define ISCSI_SUN_ISID_2    0x00    /* Sec B */
819 #define ISCSI_SUN_ISID_3    0x2A    /* Sec C - 42 = Sun's EN */
820 /*
821  * defines 4-5 are the reserved values.  These reserved values
822  * are used as the ISID for an initiator-port in MP-API and used
823  * for the send targets discovery sessions.  Byte 5 is overridden
824  * for full feature sessions.  The default values of byte 5 for a
825  * full feature session is 0.  When MS/T is enabled with more than
826  * one session this byte 5 will increment > 0 up to
827  * ISCSI_MAX_CONFIG_SESSIONS.
828  */
829 #define ISCSI_SUN_ISID_4    0x00
830 #define ISCSI_SUN_ISID_5    0xFF

832 #define ISCSI_DEFAULT_SESS_BOUND        B_FALSE
833 #define ISCSI_DEFAULT_SESS_NUM          1

835 typedef enum iscsi_enum_status {
836         ISCSI_SESS_ENUM_FREE            =       0,
837         ISCSI_SESS_ENUM_INPROG,
838         ISCSI_SESS_ENUM_DONE
839 } iscsi_enum_status_t;

841 typedef enum iscsi_enum_result {
842         ISCSI_SESS_ENUM_COMPLETE        =       0,
843         ISCSI_SESS_ENUM_PARTIAL,
844         ISCSI_SESS_ENUM_IOFAIL,
845         ISCSI_SESS_ENUM_SUBMITTED,
846         ISCSI_SESS_ENUM_SUBFAIL,
847         ISCSI_SESS_ENUM_GONE,
848         ISCSI_SESS_ENUM_TUR_FAIL
849 } iscsi_enum_result_t;

851 /*
852  * iSCSI Session(Target) Structure
853  */
```

```
854 typedef struct iscsi_sess {
855         uint32_t                sess_sig;

857         iscsi_sess_state_t      sess_state;
858         iscsi_sess_state_t      sess_prev_state;
859         clock_t                 sess_state_lbolt;
860         /* protects the session state and synchronizes the state machine */
861         krwlock_t               sess_state_rwlock;

863         /*
864          * Associated target OID.
865          */
866         uint32_t                sess_target_oid;

868         /*
869          * Session OID.  Used by IMA, interfaces and exported as
870          * TARGET_PROP which is checked by the NDI.  In addition
871          * this is used in our tran_lun_init function.
872          */
873         uint32_t                sess_oid;

875         struct iscsi_sess       *sess_next;
876         struct iscsi_hba        *sess_hba;

878         /* list of all luns relating to session */
879         struct iscsi_lun        *sess_lun_list;
880         krwlock_t               sess_lun_list_rwlock;

882         /* list of all connections relating to session */
883         struct iscsi_conn       *sess_conn_list;
884         struct iscsi_conn       *sess_conn_list_last_ptr;
885         /* pointer to active connection in session */
886         struct iscsi_conn       *sess_conn_act;
887         krwlock_t               sess_conn_list_rwlock;

889         /* Connection ID for next connection to be added to session */
890         uint32_t                sess_conn_next_cid;

892         /*
893          * last time any connection on this session received
894          * data from the target.
895          */
896         clock_t                 sess_rx_lbolt;

898         clock_t                 sess_failure_lbolt;

900         int                     sess_storm_delay;

902         /*
903          * sess_cmdsn_mutex protects the cmdsn and itt table/values
904          * Cmdsn isn't that big of a problem yet since we only have
905          * one connection but in the future we will need to ensure
906          * this locking is working so keep the sequence numbers in
907          * sync on the wire.
908          *
909          * We also use this lock to protect the ITT table and it's
910          * values.  We need to make sure someone doesn't assign
911          * a duplicate ITT value or cell to a command.  Also we
912          * need to make sure when someone is looking up an ITT
913          * that the command is still in that correct queue location.
914          */
915         kmutex_t                sess_cmdsn_mutex;

917         /*
918          * iSCSI command sequencing / windowing.  The next
919          * command to be sent via the pending queue will
```

```
 920             * get the sess_cmdsn.  If the maxcmdsn is less
 921             * than the next cmdsn then the iSCSI window is
 922             * closed and this command cannot be sent yet.
 923             * Most iscsi cmd responses from the target carry
 924             * a new maxcmdsn.  If this new maxcmdsn is greater
 925             * than the sess_maxcmdsn we will update it's value
 926             * and set a timer to fire in one tick and reprocess
 927             * the pending queue.
 928             *
 929             * The expcmdsn.   Is the value the target expects
 930             * to be sent for my next cmdsn.  If the expcmdsn
 931             * and the cmdsn get out of sync this could denote
 932             * a communication problem.
 933             */
 934            uint32_t                sess_cmdsn;
 935            uint32_t                sess_expcmdsn;
 936            uint32_t                sess_maxcmdsn;

 938            /* Next Initiator Task Tag (ITT) to use */
 939            uint32_t                sess_itt;
 940            /*
 941             * The session iscsi_cmd table is used to a fast performance
 942             * lookup of an ITT to a iscsi_cmd when we receive an iSCSI
 943             * PDU from the wire.  To reserve a location in the sess_cmd_table
 944             * we try the sess_itt % ISCSI_CMD_TABLE_SIZE if this cmd table
 945             * cell is already full.  Then increament the sess_itt and
 946             * try to get the cell position again, repeat until an empty
 947             * cell is found.  Once an empty cell is found place your
 948             * scsi_cmd point into the cell to reserve the location.  This
 949             * selection process should be done while holding the session's
 950             * mutex.
 951             */
 952            struct iscsi_cmd        *sess_cmd_table[ISCSI_CMD_TABLE_SIZE];
 953            int                     sess_cmd_table_count;

 955            /*
 956             * The pending queue contains all iscsi_cmds that require an
 957             * open MaxCmdSn window to be put on the wire and haven't
 958             * been placed on the wire.  Once placed on the wire they
 959             * will be moved to a connections specific active queue.
 960             */
 961            iscsi_queue_t           sess_queue_pending;

 963            iscsi_error_t           sess_last_err;

 965            iscsi_queue_t           sess_queue_completion;
 966            /* configured login parameters */
 967            iscsi_login_params_t    sess_params;

 969            /* general iSCSI protocol/session info */
 970            uchar_t                 sess_name[ISCSI_MAX_NAME_LEN];
 971            int                     sess_name_length;
 972            char                    sess_alias[ISCSI_MAX_NAME_LEN];
 973            int                     sess_alias_length;
 974            iSCSIDiscoveryMethod_t  sess_discovered_by;
 975            iscsi_sockaddr_t        sess_discovered_addr;
 976            uchar_t                 sess_isid[ISCSI_ISID_LEN]; /* Session ID */
 977            uint16_t                sess_tsid; /* Target ID */
 978            /*
 979             * If the target portal group tag(TPGT) is equal to ISCSI_DEFAULT_TPGT
 980             * then the initiator will accept a successful login with any TPGT
 981             * specified by the target.  If a none default TPGT is configured
 982             * then we will only successfully accept a login with that matching
 983             * TPGT value.
 984             */
 985            int                     sess_tpgt_conf;
```

```
 986            /* This field records the negotiated TPGT value, preserved for dtrace */
 987            int                     sess_tpgt_nego;

 989            /*
 990             * Authentication information.
 991             *
 992             * DCW: Again IMA seems to take a session view at this
 993             * information.
 994             */
 995            iscsi_auth_t            sess_auth;

 997            /* Statistics */
 998            struct {
 999                    kstat_t                 *ks;
1000                    iscsi_sess_stats_t      ks_data;
1001                    kstat_t                 *ks_io;
1002                    kstat_io_t              ks_io_data;
1003                    kmutex_t                ks_io_lock;
1004            } stats;

1006            iscsi_thread_t          *sess_ic_thread;
1007            boolean_t               sess_window_open;
1008            boolean_t               sess_boot;
1009            iscsi_sess_type_t       sess_type;

1011            ddi_taskq_t             *sess_login_taskq;

1013            iscsi_thread_t          *sess_wd_thread;

1015            sm_audit_buf_t          sess_state_audit;

1017            kmutex_t                sess_reset_mutex;

1019            boolean_t               sess_reset_in_progress;

1021            boolean_t               sess_boot_nic_reset;
1022            kmutex_t                sess_enum_lock;
1023            kcondvar_t              sess_enum_cv;
1024            iscsi_enum_status_t     sess_enum_status;
1025            iscsi_enum_result_t     sess_enum_result;
1026            uint32_t                sess_enum_result_count;
1027            ddi_taskq_t             *sess_enum_taskq;

1029            kmutex_t                sess_state_wmutex;
1030            kcondvar_t              sess_state_wcv;
1031            boolean_t               sess_state_hasw;

1033            /* to accelerate the state change in case of new event */
1034            volatile uint32_t       sess_state_event_count;
1035 } iscsi_sess_t;

1037 /*
1038  * This structure will be used to store sessions to be online
1039  * during normal login operation.
1040  */
1041 typedef struct iscsi_sess_list {
1042         iscsi_sess_t            *session;
1043         struct iscsi_sess_list  *next;
1044 } iscsi_sess_list_t;

1046 /*
1047  * iSCSI client notify task context for deferred IDM notifications processing
1048  */
1049 typedef struct iscsi_cn_task {
1050         idm_conn_t              *ct_ic;
1051         idm_client_notify_t     ct_icn;
```

```
1052            uintptr_t               ct_data;
1053 } iscsi_cn_task_t;

1055 /*
1056  * iscsi_network
1057  */
1058 typedef struct iscsi_network {
1059         void* (*socket)(int domain, int, int);
1060         int (*bind)(void *, struct sockaddr *, int, int, int);
1061         int (*connect)(void *, struct sockaddr *, int, int, int);
1062         int (*listen)(void *, int);
1063         void* (*accept)(void *, struct sockaddr *, int *);
1064         int (*getsockname)(void *, struct sockaddr *, socklen_t *);
1065         int (*getsockopt)(void *, int, int, void *, int *, int);
1066         int (*setsockopt)(void *, int, int, void *, int);
1067         int (*shutdown)(void *, int);
1068         void (*close)(void *);

1070         size_t (*poll)(void *, clock_t);
1071         size_t (*sendmsg)(void *, struct msghdr *);
1072         size_t (*recvmsg)(void *, struct msghdr *, int);

1074         iscsi_status_t (*sendpdu)(void *, iscsi_hdr_t *, char *, int);
1075         iscsi_status_t (*recvdata)(void *, iscsi_hdr_t *, char *,
1076             int, int, int);
1077         iscsi_status_t (*recvhdr)(void *, iscsi_hdr_t *, int, int, int);

1079         struct {
1080                 int                     sndbuf;
1081                 int                     rcvbuf;
1082                 int                     nodelay;
1083                 int                     conn_notify_threshold;
1084                 int                     conn_abort_threshold;
1085                 int                     abort_threshold;
1086         } tweaks;
1087 } iscsi_network_t;

1089 #define ISCSI_NET_HEADER_DIGEST 0x00000001
1090 #define ISCSI_NET_DATA_DIGEST   0x00000002

1092 extern iscsi_network_t *iscsi_net;

1094 /*
1095  * If we get bus_config requests in less than 5 seconds
1096  * apart skip the name services re-discovery and just
1097  * complete the requested logins.  This protects against
1098  * bus_config storms from stale /dev links.
1099  */
1100 #define ISCSI_CONFIG_STORM_DELAY_DEFAULT     5

1102 /*
1103  * iSCSI HBA Structure
1104  */
1105 typedef struct iscsi_hba {
1106         uint32_t                hba_sig;
1107         dev_info_t              *hba_dip;       /* dev info ptr */
1108         scsi_hba_tran_t         *hba_tran;      /* scsi tran ptr */
1109         ldi_ident_t             hba_li;

1111         struct iscsi_sess       *hba_sess_list; /* sess. list for hba */
1112         krwlock_t               hba_sess_list_rwlock; /* protect sess. list */

1114         /* lbolt of the last time we received a config request */
1115         clock_t                 hba_config_lbolt;
1116         /* current number of seconds to protect against bus config storms */
1117         int                     hba_config_storm_delay;
```

```
1119         /* general iSCSI protocol hba/initiator info */
1120         uchar_t                 hba_name[ISCSI_MAX_NAME_LEN];
1121         int                     hba_name_length;
1122         uchar_t                 hba_alias[ISCSI_MAX_NAME_LEN];
1123         int                     hba_alias_length;

1125         /* Default SessionID for HBA */
1126         uchar_t                 hba_isid[ISCSI_ISID_LEN];

1128         /* Default HBA wide settings */
1129         iscsi_login_params_t    hba_params;

1131         /*
1132          * There's only one HBA and it's set to ISCSI_INITIATOR_OID
1133          * (value of 1) at the beginning of time.
1134          */
1135         uint32_t                hba_oid;

1137         /*
1138          * Keep track of which events have been sent. User daemons request
1139          * this information so they don't wait for events which they won't
1140          * see.
1141          */
1142         kmutex_t                hba_discovery_events_mutex;
1143         iSCSIDiscoveryMethod_t  hba_discovery_events;
1144         boolean_t               hba_discovery_in_progress;

1146         boolean_t               hba_mpxio_enabled; /* mpxio-enabled */
1147         /* if the persistent store is loaded */
1148         boolean_t               hba_persistent_loaded;

1150         /*
1151          * Ensures only one SendTargets operation occurs at a time
1152          */
1153         ksema_t                 hba_sendtgts_semaphore;

1155         /*
1156          * Statistics
1157          */
1158         struct {
1159                 kstat_t                 *ks;
1160                 iscsi_hba_stats_t       ks_data;
1161         } stats;

1163         /*
1164          * track/control the service status and client
1165          *
1166          * service- service online ensures the operational of cli
1167          *        - and the availability of iSCSI discovery/devices
1168          *        - so obviously offline means the unusable of cli
1169          *        - , disabling of all discovery methods and to offline
1170          *        - all discovered devices
1171          *
1172          * client - here the client actually means 'exclusive client'
1173          *        - for operations these clients take may conflict
1174          *        - with the changing of service status and therefore
1175          *        - need to be exclusive
1176          *
1177          * The service has three status:
1178          *      ISCSI_SERVICE_ENABLED   -       client is permitted to
1179          *      -                               request service
1180          *
1181          *      ISCSI_SERVICE_DISABLED  -       client is not permitted to
1182          *      -                               request service
1183          *
```

```
1184             *       ISCSI_SERVICE_TRANSITION -      client must wait for
1185             *                               -       one of above two statuses
1186             *
1187             * The hba_service_client_count tracks the number of
1188             * current clients, it increases with new clients and decreases
1189             * with leaving clients. It stops to increase once the
1190             * ISCSI_SERVICE_TRANSITION is set, and causes later clients be
1191             * blocked there.
1192             *
1193             * The status of the service can only be changed when the number
1194             * of current clients reaches zero.
1195             *
1196             * Clients include:
1197             *       iscsi_ioctl
1198             *       iscsi_tran_bus_config
1199             *       iscsi_tran_bus_unconfig
1200             *       isns_scn_callback
1201             */
1202            kmutex_t                hba_service_lock;
1203            kcondvar_t              hba_service_cv;
1204            uint32_t                hba_service_status;
1205            uint32_t                hba_service_client_count;

1207            /* Default HBA tunable settings */
1208            iscsi_tunable_params_t  hba_tunable_params;
1209            boolean_t               hba_service_status_overwrite;
1210 } iscsi_hba_t;

1212 /*
1213  * +------------------------------------------------------------------+
1214  * | iSCSI prototypes                                                 |
1215  * +------------------------------------------------------------------+
1216  */

1218 /* IDM client callback entry points */
1219 idm_rx_pdu_cb_t iscsi_rx_scsi_rsp;
1220 idm_rx_pdu_cb_t iscsi_rx_misc_pdu;
1221 idm_rx_pdu_error_cb_t iscsi_rx_error_pdu;
1222 idm_build_hdr_cb_t iscsi_build_hdr;
1223 idm_task_cb_t iscsi_task_aborted;
1224 idm_client_notify_cb_t iscsi_client_notify;

1226 /* iscsi_io.c */
1227 int iscsi_sna_lte(uint32_t n1, uint32_t n2);
1228 char *iscsi_get_next_text(char *data, int data_length, char *curr_text);

1230 void iscsi_ic_thread(iscsi_thread_t *thread, void *arg);
1231 void iscsi_tx_thread(iscsi_thread_t *thread, void *arg);
1232 void iscsi_wd_thread(iscsi_thread_t *thread, void *arg);

1234 iscsi_status_t iscsi_tx_cmd(iscsi_sess_t *isp, iscsi_cmd_t *icmdp);

1236 void iscsi_task_cleanup(int opcode, iscsi_cmd_t *icmdp);

1238 void iscsi_handle_abort(void *arg);
1239 iscsi_status_t iscsi_handle_reset(iscsi_sess_t *isp, int level,
1240     iscsi_lun_t *ilp);
1241 iscsi_status_t iscsi_handle_logout(iscsi_conn_t *icp);
1242 iscsi_status_t iscsi_handle_passthru(iscsi_sess_t *isp, uint16_t lun,
1243     struct uscsi_cmd *ucmdp);
1244 iscsi_status_t iscsi_handle_text(iscsi_conn_t *icp,
1245     char *buf, uint32_t buf_len, uint32_t data_len, uint32_t *rx_data_len);

1247 void iscsi_iodone(iscsi_sess_t *isp, iscsi_cmd_t *icmdp);

1249 /* iscsi_crc.c */
```

```
1250 uint32_t iscsi_crc32c(void *address, unsigned long length);
1251 uint32_t iscsi_crc32c_continued(void *address, unsigned long length,
1252     uint32_t crc);

1254 /* iscsi_queue.c */
1255 void iscsi_init_queue(iscsi_queue_t *queue);
1256 void iscsi_destroy_queue(iscsi_queue_t *queue);
1257 void iscsi_enqueue_pending_cmd(iscsi_sess_t *isp, iscsi_cmd_t *icmdp);
1258 void iscsi_dequeue_pending_cmd(iscsi_sess_t *isp, iscsi_cmd_t *icmdp);
1259 void iscsi_enqueue_active_cmd(iscsi_conn_t *icp, iscsi_cmd_t *icmdp);
1260 void iscsi_dequeue_active_cmd(iscsi_conn_t *icp, iscsi_cmd_t *icmdp);
1261 void iscsi_enqueue_idm_aborting_cmd(iscsi_conn_t *icp, iscsi_cmd_t *icmdp);
1262 void iscsi_dequeue_idm_aborting_cmd(iscsi_conn_t *icp, iscsi_cmd_t *icmdp);
1263 void iscsi_enqueue_completed_cmd(iscsi_sess_t *isp, iscsi_cmd_t *icmdp);
1264 iscsi_status_t iscsi_dequeue_cmd(iscsi_cmd_t **, iscsi_cmd_t **, iscsi_cmd_t *);
1265 void iscsi_move_queue(iscsi_queue_t *src_queue, iscsi_queue_t *dst_queue);
1266 void iscsi_enqueue_cmd_head(iscsi_cmd_t **, iscsi_cmd_t **,
1267     iscsi_cmd_t *);

1269 /* iscsi_login.c */
1270 iscsi_status_t iscsi_login_start(void *arg);
1271 void iscsi_login_update_state(iscsi_conn_t *icp,
1272     iscsi_login_state_t next_state);
1273 void iscsi_login_update_state_locked(iscsi_conn_t *icp,
1274     iscsi_login_state_t next_state);


1277 /* iscsi_stats.c */
1278 boolean_t iscsi_hba_kstat_init(struct iscsi_hba *ihp);
1279 boolean_t iscsi_hba_kstat_term(struct iscsi_hba *ihp);
1280 boolean_t iscsi_sess_kstat_init(struct iscsi_sess *isp);
1281 boolean_t iscsi_sess_kstat_term(struct iscsi_sess *isp);
1282 boolean_t iscsi_conn_kstat_init(struct iscsi_conn      *icp);
1283 void iscsi_conn_kstat_term(struct iscsi_conn *icp);

1285 /* iscsi_net.c */
1286 void iscsi_net_init();
1287 void iscsi_net_fini();
1288 iscsi_status_t iscsi_net_interface(boolean_t reset);

1290 /* iscsi_sess.c */
1291 iscsi_sess_t *iscsi_sess_create(iscsi_hba_t *ihp,
1292     iSCSIDiscoveryMethod_t method, struct sockaddr *addr_dsc,
1293     char *target_name, int tpgt, uchar_t isid_lsb,
1294     iscsi_sess_type_t type, uint32_t *oid);
1295 void iscsi_sess_online(void *arg);
1296 int iscsi_sess_get(uint32_t oid, iscsi_hba_t *ihp, iscsi_sess_t **ispp);
1297 iscsi_status_t iscsi_sess_destroy(iscsi_sess_t *isp);
1298 void iscsi_sess_state_machine(iscsi_sess_t *isp, iscsi_sess_event_t event,
1299     uint32_t event_count);
1300 char *iscsi_sess_state_str(iscsi_sess_state_t state);
1301 boolean_t iscsi_sess_set_auth(iscsi_sess_t *isp);
1302 iscsi_status_t iscsi_sess_reserve_scsi_itt(iscsi_cmd_t *icmdp);
1303 void iscsi_sess_release_scsi_itt(iscsi_cmd_t *icmdp);
1304 iscsi_status_t iscsi_sess_reserve_itt(iscsi_sess_t *isp, iscsi_cmd_t *icmdp);
1305 void iscsi_sess_release_itt(iscsi_sess_t *isp, iscsi_cmd_t *icmdp);
1306 void iscsi_sess_redrive_io(iscsi_sess_t *isp);
1307 int iscsi_sess_get_by_target(uint32_t target_oid, iscsi_hba_t *ihp,
1308     iscsi_sess_t **ispp);
1309 iscsi_enum_result_t iscsi_sess_enum_request(iscsi_sess_t *isp,
1310     boolean_t wait, uint32_t event_count);
1311 iscsi_enum_result_t iscsi_sess_enum_query(iscsi_sess_t *isp);
1312 void iscsi_sess_enter_state_zone(iscsi_sess_t *isp);
1313 void iscsi_sess_exit_state_zone(iscsi_sess_t *isp);


1315 /* iscsi_conn.c */
```

```
1316 iscsi_status_t iscsi_conn_create(struct sockaddr *addr, iscsi_sess_t *isp,
1317     iscsi_conn_t **icpp);
1318 iscsi_status_t iscsi_conn_online(iscsi_conn_t *icp);
1319 iscsi_status_t iscsi_conn_offline(iscsi_conn_t *icp);
1320 iscsi_status_t iscsi_conn_destroy(iscsi_conn_t *icp);
1321 void iscsi_conn_set_login_min_max(iscsi_conn_t *icp, int min, int max);
1322 iscsi_status_t iscsi_conn_sync_params(iscsi_conn_t *icp);
1323 void iscsi_conn_retry(iscsi_sess_t *isp, iscsi_conn_t *icp);
1324 void iscsi_conn_update_state(iscsi_conn_t *icp, iscsi_conn_state_t next_state);
1325 void iscsi_conn_update_state_locked(iscsi_conn_t *icp,
1326                     iscsi_conn_state_t next_state);

1328 /* iscsi_lun.c */
1329 iscsi_status_t iscsi_lun_create(iscsi_sess_t *isp, uint16_t lun_num,
1330     uint8_t lun_addr_type, struct scsi_inquiry *inq, char *guid);
1331 iscsi_status_t iscsi_lun_destroy(iscsi_hba_t *ihp,
1332     iscsi_lun_t *ilp);
1333 void iscsi_lun_online(iscsi_hba_t *ihp,
1334     iscsi_lun_t *ilp);
1335 iscsi_status_t iscsi_lun_offline(iscsi_hba_t *ihp,
1336     iscsi_lun_t *ilp, boolean_t lun_free);

1338 /* iscsi_cmd.c */
1339 void iscsi_cmd_state_machine(iscsi_cmd_t *icmdp,
1340     iscsi_cmd_event_t event, void *arg);
1341 iscsi_cmd_t     *iscsi_cmd_alloc(iscsi_conn_t *icp, int km_flags);
1342 void           iscsi_cmd_free(iscsi_cmd_t *icmdp);

1344 /* iscsi_ioctl.c */
1345 void * iscsi_ioctl_copyin(caddr_t arg, int mode, size_t size);
1346 int iscsi_ioctl_copyout(void *data, size_t size, caddr_t arg, int mode);
1347 iscsi_conn_list_t *iscsi_ioctl_conn_oid_list_get_copyin(caddr_t, int);
1348 int iscsi_ioctl_conn_oid_list_get_copyout(iscsi_conn_list_t *, caddr_t, int);
1349 boolean_t iscsi_ioctl_conn_oid_list_get(iscsi_hba_t *ihp,
1350     iscsi_conn_list_t *cl);
1351 boolean_t iscsi_ioctl_conn_props_get(iscsi_hba_t *ihp, iscsi_conn_props_t *cp);
1352 int iscsi_ioctl_sendtgts_get(iscsi_hba_t *ihp, iscsi_sendtgts_list_t *stl);
1353 int iscsi_target_prop_mod(iscsi_hba_t *, iscsi_property_t *, int cmd);
1354 int iscsi_set_params(iscsi_param_set_t *, iscsi_hba_t *, boolean_t);
1355 int iscsi_get_persisted_param(uchar_t *, iscsi_param_get_t *,
1356     iscsi_login_params_t *);
1357 void iscsi_set_default_login_params(iscsi_login_params_t *params);
1358 int iscsi_ioctl_get_config_sess(iscsi_hba_t *ihp,
1359     iscsi_config_sess_t *ics);
1360 int iscsi_ioctl_set_config_sess(iscsi_hba_t *ihp,
1361     iscsi_config_sess_t *ics);
1362 int iscsi_ioctl_set_tunable_param(iscsi_hba_t *ihp,
1363     iscsi_tunable_object_t *tpss);
1364 /* ioctls  prototypes */
1365 int iscsi_get_param(iscsi_login_params_t *params,
1366     boolean_t valid_flag,
1367     iscsi_param_get_t *ipgp);

1369 /* iscsid.c */
1370 boolean_t iscsid_init(iscsi_hba_t *ihp);
1371 boolean_t iscsid_start(iscsi_hba_t *ihp);
1372 boolean_t iscsid_stop(iscsi_hba_t *ihp);
1373 void iscsid_fini();
1374 void iscsid_props(iSCSIDiscoveryProperties_t *props);
1375 boolean_t iscsid_enable_discovery(iscsi_hba_t *ihp,
1376     iSCSIDiscoveryMethod_t idm, boolean_t poke);
1377 boolean_t iscsid_disable_discovery(iscsi_hba_t *ihp,
1378     iSCSIDiscoveryMethod_t idm);
1379 void iscsid_poke_discovery(iscsi_hba_t *ihp, iSCSIDiscoveryMethod_t method);
1380 void iscsid_do_sendtgts(entry_t *discovery_addr);
1381 void iscsid_do_isns_query_one_server(
```

```
1382     iscsi_hba_t *ihp, entry_t *isns_addr);
1383 void iscsid_do_isns_query(iscsi_hba_t *ihp);
1384 void iscsid_config_one(iscsi_hba_t *ihp,
1385     char *name, boolean_t protect);
1386 void iscsid_config_all(iscsi_hba_t *ihp, boolean_t protect);
1387 void iscsid_unconfig_one(iscsi_hba_t *ihp, char *name);
1388 void iscsid_unconfig_all(iscsi_hba_t *ihp);
1389 void isns_scn_callback(void *arg);
1390 boolean_t iscsid_del(iscsi_hba_t *ihp, char *target_name,
1391     iSCSIDiscoveryMethod_t method, struct sockaddr *addr_dsc);
1392 boolean_t iscsid_login_tgt(iscsi_hba_t *ihp, char *target_name,
1393     iSCSIDiscoveryMethod_t method, struct sockaddr *addr_dsc);
1394 void iscsid_addr_to_sockaddr(int src_insize, void *src_addr, int src_port,
1395     struct sockaddr *dst_addr);
1396 void iscsid_set_default_initiator_node_settings(iscsi_hba_t *ihp,
1397     boolean_t minimal);

1399 void iscsi_send_sysevent(iscsi_hba_t *ihp, char *eventcalss,
1400     char *subclass, nvlist_t *np);
1401 boolean_t iscsi_reconfig_boot_sess(iscsi_hba_t *ihp);
1402 boolean_t iscsi_chk_bootlun_mpxio(iscsi_hba_t *ihp);
1403 boolean_t iscsi_cmp_boot_ini_name(char *name);
1404 boolean_t iscsi_cmp_boot_tgt_name(char *name);
1405 boolean_t iscsi_client_request_service(iscsi_hba_t *ihp);
1406 void iscsi_client_release_service(iscsi_hba_t *ihp);

1408 extern void bcopy(const void *s1, void *s2, size_t n);
1409 extern void bzero(void *s, size_t n);
  24 /*
  25  * Here we need a contract for inet_ntop() and inet_pton()
  26  * in uts/common/inet/ip/inet_ntop.c
  27  */
  28 extern char *inet_ntop(int af, const void *addr, char *buf, int addrlen);
  29 extern int inet_pton(int af, char *inp, void *outp);

1411 #ifdef __cplusplus
1412 }
_____unchanged_portion_omitted_
```

```
**********************************************************
    9400 Wed Sep 26 12:51:50 2012
new/usr/src/uts/common/io/scsi/adapters/iscsi/iscsi_doorclt.c
inet_pton
**********************************************************
   1 /*
   2  * CDDL HEADER START
   3  *
   4  * The contents of this file are subject to the terms of the
   5  * Common Development and Distribution License (the "License").
   6  * You may not use this file except in compliance with the License.
   7  *
   8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
   9  * or http://www.opensolaris.org/os/licensing.
  10  * See the License for the specific language governing permissions
  11  * and limitations under the License.
  12  *
  13  * When distributing Covered Code, include this CDDL HEADER in each
  14  * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
  15  * If applicable, add the following below this CDDL HEADER, with the
  16  * fields enclosed by brackets "[]" replaced with your own identifying
  17  * information: Portions Copyright [yyyy] [name of copyright owner]
  18  *
  19  * CDDL HEADER END
  20  */
  21 /*
  22  * Copyright 2000 by Cisco Systems, Inc.  All rights reserved.
  23  * Copyright 2009 Sun Microsystems, Inc.  All rights reserved.
  24  * Use is subject to license terms.
  25  * Copyright 2012 Nexenta Systems, Inc. All rights reserved.
  26 #endif /* ! codereview */
  27  *
  28  * iSCSI Software Initiator
  29  */

  31 #include <sys/types.h>
  32 #include <sys/errno.h>
  33 #include <sys/conf.h>
  34 #include <sys/cmn_err.h>
  35 #include <sys/stat.h>
  36 #include <sys/pathname.h>
  37 #include <sys/door.h>
  38 #include <sys/kmem.h>
  39 #include <sys/socket.h>
  40 #include <sys/fs/snode.h>
  41 #include <netinet/in.h>

  43 #include <sys/scsi/adapters/iscsi_door.h>
  44 #include "iscsi.h"

  46 #define ISCSI_DOOR_MAX_SEMA_VALUE        16

  48 static boolean_t        iscsi_door_init = B_FALSE;
  49 static ksema_t          iscsi_door_sema;
  50 static krwlock_t        iscsi_door_lock;
  51 static door_handle_t    iscsi_door_handle;

  53 typedef struct _mybuffer {
  54         size_t          signature;
  55         size_t          size;
  56 } mybuffer_t;

  58 /*
  59  * iscsi_door_ini
  60  *
  61  * This function initializes the variables needed to handle the door upcall.
```

```
  62  */
  63 boolean_t
  64 iscsi_door_ini(void)
  65 {
  66         ASSERT(!iscsi_door_init);
  67         if (!iscsi_door_init) {
  68                 rw_init(
  69                     &iscsi_door_lock,
  70                     NULL,
  71                     RW_DRIVER,
  72                     NULL);

  74                 sema_init(
  75                     &iscsi_door_sema,
  76                     ISCSI_DOOR_MAX_SEMA_VALUE,
  77                     NULL,
  78                     SEMA_DRIVER,
  79                     NULL);

  81                 iscsi_door_handle = NULL;
  82                 iscsi_door_init = B_TRUE;
  83                 return (B_TRUE);
  84         }
  85         return (B_FALSE);
  86 }

  88 /*
  89  * iscsi_door_term
  90  *
  91  * This function releases the resources allocated to handle the door
  92  * upcall.  It disconnects from the door if currently connected.
  93  */
  94 boolean_t
  95 iscsi_door_term(void)
  96 {
  97         ASSERT(iscsi_door_init);
  98         if (iscsi_door_init) {
  99                 iscsi_door_init = B_FALSE;
 100                 iscsi_door_unbind();
 101                 rw_destroy(&iscsi_door_lock);
 102                 sema_destroy(&iscsi_door_sema);
 103                 return (B_TRUE);
 104         }
 105         return (B_FALSE);
 106 }

 108 /*
 109  * iscsi_door_bind
 110  *
 111  * This function tries to connect the iscsi_door.  If it succeeds
 112  * it keeps the vnode.
 113  */
 114 boolean_t
 115 iscsi_door_bind(
 116         int             did
 117 )
 118 {
 119         door_handle_t   new_handle;

 121         new_handle = door_ki_lookup(did);
 122         if (new_handle == NULL) {
 123                 /* The lookup failed. */
 124                 return (B_FALSE);
 125         }

 127         /* The new handle is stored.  If we had one, it is released. */
```

```
 128            rw_enter(&iscsi_door_lock, RW_WRITER);
 129            if (iscsi_door_handle != NULL) {
 130                    door_ki_rele(iscsi_door_handle);
 131            }
 132            iscsi_door_handle = new_handle;
 133            rw_exit(&iscsi_door_lock);

 135            return (B_TRUE);
 136 }

 138 /*
 139  * iscsi_door_unbind
 140  *
 141  * This function releases the current door handle.
 142  */
 143 void
 144 iscsi_door_unbind(void)
 145 {
 146            rw_enter(&iscsi_door_lock, RW_WRITER);
 147            if (iscsi_door_handle != NULL) {
 148                    door_ki_rele(iscsi_door_handle);
 149                    iscsi_door_handle = NULL;
 150            }
 151            rw_exit(&iscsi_door_lock);
 152 }

 154 /*
 155  * iscsi_door_upcall
 156  *
 157  * This function tries to call the iscsi_door.
 158  */
 159 static
 160 boolean_t
 161 iscsi_door_upcall(door_arg_t *arg)
 162 {
 163            int     error;

 165            /*
 166             * This semaphore limits the number of simultaneous calls
 167             * to the door.
 168             */
 169            sema_p(&iscsi_door_sema);
 170            /*
 171             * The mutex protecting the iscsi_door_handle is entered.
 172             */
 173            rw_enter(&iscsi_door_lock, RW_READER);

 175            if (iscsi_door_handle == NULL) {
 176                    /* There's no door handle. */
 177                    rw_exit(&iscsi_door_lock);
 178                    sema_v(&iscsi_door_sema);
 179                    return (B_FALSE);
 180            }
 181            error = door_ki_upcall(iscsi_door_handle, arg);

 183            rw_exit(&iscsi_door_lock);
 184            sema_v(&iscsi_door_sema);

 186            if (error != 0) {
 187                    return (B_FALSE);
 188            } else {
 189                    return (B_TRUE);
 190            }
 191 }

 193 /*
```

```
 194  * kfreehostent
 195  *
 196  * This function frees the memory returned by kgetipnodebyname.
 197  */
 198 void
 199 kfreehostent(
 200            struct hostent          *hptr
 201 )
 202 {
 203            mybuffer_t              *buffer;

 205            ASSERT(hptr != NULL);
 206            if (hptr) {
 207                    buffer = (mybuffer_t *)((char *)hptr - sizeof (mybuffer_t));
 208                    ASSERT(buffer->signature == ISCSI_DOOR_REQ_SIGNATURE);
 209                    if (buffer->signature == ISCSI_DOOR_REQ_SIGNATURE) {
 210                            kmem_free((void *)buffer, buffer->size);
 211                            return;
 212                    }
 213            }
 214            /* A message should be logged here. */
 215 }

 217 /*
 218  * kgetipnodebyname
 219  *
 220  * This function builds a request that will be sent to the iscsi_door.
 221  * The iSCSI door after receiving the request calls getipnodebyaddr().
 222  * for more information on the input, output parameter and return value,
 223  * consult the man page for getipnodebyname().
 224  *
 225  * Before calling the iscsi door this function tries to do the conversion
 226  * locally.  If a name resolution is needed the iscsi door is called.
 227  *
 228  * There's some limitations to the information returned by this function.
 229  * Only one address of the address list returned by getipnodebyname() is
 230  * returned.  The other parameters of the structure should be ignored.
 231  */
 232 struct hostent *
 233 kgetipnodebyname(
 234            const char      *name,
 235            int             af,
 236            int             flags,
 237            int             *error_num
 238 )
 239 {
 240            door_arg_t              arg;
 241            mybuffer_t              *buffer;
 242            size_t                  msg_size = ISCSI_DOOR_MAX_DATA_SIZE;
 243            size_t                  hostent_size = ISCSI_DOOR_MAX_DATA_SIZE;
 244            size_t                  buffer_size;
 245            getipnodebyname_req_t   *req;
 246            getipnodebyname_cnf_t   *cnf;
 247            struct hostent          *hptr;
  25            int                     i;
  26            uint16_t                *swap;


 250            buffer_size = msg_size + hostent_size + sizeof (mybuffer_t);
 251            buffer = (mybuffer_t *)kmem_zalloc(buffer_size, KM_SLEEP);

 253            if (buffer) {

 255                    /*
 256                     * The buffer was successfully allocated.
 257                     *
```

```
 258                 *          Buffer
 259                 *
 260                 * +-------------------+ <--- buffer
 261                 * |    mybuffer_t     |
 262                 * +-------------------+ <--- hptr
 263                 * |                   |
 264                 * |                   |
 265                 * |    hostent_size   |
 266                 * |                   |
 267                 * |                   |
 268                 * |                   |
 269                 * +-------------------+ <--- req, cnf
 270                 * |                   |
 271                 * |                   |
 272                 * |                   |
 273                 * |    msg_size       |
 274                 * |                   |
 275                 * |                   |
 276                 * |                   |
 277                 * +-------------------+
 278                 */
 279                buffer->signature = ISCSI_DOOR_REQ_SIGNATURE;
 280                buffer->size = buffer_size;

 282                hptr = (struct hostent *)((char *)buffer + sizeof (mybuffer_t));
 283                req = (getipnodebyname_req_t *)((char *)hptr + hostent_size);
 284                cnf = (getipnodebyname_cnf_t *)((char *)hptr + hostent_size);

 286                hostent_size -= sizeof (struct hostent);

 288                /*
 289                 * We try first locally.  If the conversion cannot be done
 290                 * by inet_pton the door is called.
 291                 * The cnf address is used as output buffer.
 292                 * inet_pton returns '1' if the conversion was successful.
 293                 */
 294                switch (af) {
 295                case AF_INET:
 296                        hptr->h_length = sizeof (struct in_addr);
 297                        break;
 298                case AF_INET6:
 299                        hptr->h_length = sizeof (struct in6_addr);
 300                        break;
 301                default:
 302                        kfreehostent(hptr);
 303                        *error_num = NO_RECOVERY;
 304                        return (NULL);
 305                }
 306                if ((msg_size < hptr->h_length) ||
 307                    (hostent_size < sizeof (char *))) {
 308                        kfreehostent(hptr);
 309                        *error_num = NO_RECOVERY;
 310                        return (NULL);
 311                }
 312                if (inet_pton(af, (char *)name, cnf) == 1) {
 313                        /*
 314                         * inet_pton converted the string successfully.
  94                         * reset to network order.  swaps based on nfs code
 315                         */
  96                        if (af == AF_INET) {
  97                                *((uint32_t *)cnf) = htonl(*((uint32_t *)cnf));
  98                        } else {
  99                                for (swap = ((void *)cnf), i = 0;
 100                                    i < hptr->h_length / sizeof (uint16_t);
 101                                    i++) {
 102                                        swap[i] = htons(swap[i]);
```

```
 103                                }
 104                        }
 316                        hptr->h_addrtype = af;
 317                        hptr->h_addr_list = (char **)((char *)hptr +
 318                            sizeof (struct hostent));
 319                        *hptr->h_addr_list = (char *)cnf;
 320                        return (hptr);
 321                }

 323                /*
 324                 * The name couldn't ne converted by inet_pton.  The door is
 325                 * called.
 326                 */

 328                /* Header initialization. */
 329                req->hdr.signature = ISCSI_DOOR_REQ_SIGNATURE;
 330                req->hdr.version = ISCSI_DOOR_REQ_VERSION_1;
 331                req->hdr.opcode = ISCSI_DOOR_GETIPNODEBYNAME_REQ;

 333                /* Body initialization. */
 334                req->name_length = strlen(name);
 335                if (req->name_length >
 336                    (msg_size - sizeof (getipnodebyname_req_t) - 1)) {
 337                        kfreehostent(hptr);
 338                        *error_num = NO_RECOVERY;
 339                        return (NULL);
 340                }

 342                req->name_offset = sizeof (getipnodebyname_req_t);
 343                req->af = af;
 344                req->flags = flags;
 345                bcopy(
 346                    name,
 347                    ((char *)req + req->name_offset),
 348                    req->name_length);

 350                /* Door argument initialization. */
 351                arg.data_ptr = (char *)req;
 352                arg.data_size = msg_size;
 353                arg.desc_num = 0;
 354                arg.desc_ptr = NULL;
 355                arg.rbuf = (char *)cnf;
 356                arg.rsize = msg_size;

 358                if (iscsi_door_upcall(&arg) == B_FALSE) {
 359                        /* The door call failed */
 360                        kfreehostent(hptr);
 361                        *error_num = NO_RECOVERY;
 362                        return (NULL);
 363                }

 365                /*
 366                 * The door call itself was successful.  The value returned
 367                 * in arg.rbuf should be cnf, but we never know.
 368                 */
 369                cnf = (getipnodebyname_cnf_t *)arg.rbuf;

 371                if ((cnf == NULL) ||
 372                    (arg.rsize < sizeof (getipnodebyname_cnf_t)) ||
 373                    (cnf->hdr.signature != ISCSI_DOOR_REQ_SIGNATURE) ||
 374                    (cnf->hdr.version != ISCSI_DOOR_REQ_VERSION_1) ||
 375                    (cnf->hdr.opcode != ISCSI_DOOR_GETIPNODEBYNAME_CNF) ||
 376                    ((cnf->hdr.status != ISCSI_DOOR_STATUS_SUCCESS) &&
 377                    (cnf->hdr.status != ISCSI_DOOR_STATUS_MORE))) {
 378                        /* The door didn't like the request */
 379                        kfreehostent(hptr);
```

```
 380                                *error_num = NO_RECOVERY;
 381                                return (NULL);
 382                        }

 384                        if (cnf->h_addr_list_length == 0) {
 385                                kfreehostent(hptr);
 386                                *error_num = HOST_NOT_FOUND;
 387                                return (NULL);
 388                        }

 390                        hptr->h_addrtype = cnf->h_addrtype;
 391                        hptr->h_length = cnf->h_addrlen;
 392                        hptr->h_addr_list = (char **)((char *)hptr +
 393                            sizeof (struct hostent));
 394                        *hptr->h_addr_list = ((char *)cnf + cnf->h_addr_list_offset);
 395                        return (hptr);
 396                } else {
 397                        *error_num = NO_RECOVERY;
 398                        return (NULL);
 399                }
 400 }
_____unchanged_portion_omitted_
```