```
**********************************************************
   57931 Mon May  5 11:11:31 2014
new/usr/src/uts/common/io/mac/mac_protect.c
4788 mac shouldn't abuse ddi_get_time(9f)
**********************************************************
   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) 2010, Oracle and/or its affiliates. All rights reserved.
  24  */
  25 /*
  26  * Copyright 2014 Nexenta Systems, Inc.  All rights reserved.
  27  */
  28 #endif /* ! codereview */

  30 #include <sys/strsun.h>
  31 #include <sys/sdt.h>
  32 #include <sys/mac.h>
  33 #include <sys/mac_impl.h>
  34 #include <sys/mac_client_impl.h>
  35 #include <sys/mac_client_priv.h>
  36 #include <sys/ethernet.h>
  37 #include <sys/vlan.h>
  38 #include <sys/dlpi.h>
  39 #include <sys/avl.h>
  40 #include <inet/ip.h>
  41 #include <inet/ip6.h>
  42 #include <inet/arp.h>
  43 #include <netinet/arp.h>
  44 #include <netinet/udp.h>
  45 #include <netinet/dhcp.h>
  46 #include <netinet/dhcp6.h>

  48 /*
  49  * Implementation overview for DHCP address detection
  50  *
  51  * The purpose of DHCP address detection is to relieve the user of having to
  52  * manually configure static IP addresses when ip-nospoof protection is turned
  53  * on. To achieve this, the mac layer needs to intercept DHCP packets to
  54  * determine the assigned IP addresses.
  55  *
  56  * A DHCP handshake between client and server typically requires at least
  57  * 4 messages:
  58  *
  59  * 1. DISCOVER - client attempts to locate DHCP servers via a
  60  *                 broadcast message to its subnet.
  61  * 2. OFFER    - server responds to client with an IP address and
```

```
  62  *                 other parameters.
  63  * 3. REQUEST  - client requests the offered address.
  64  * 4. ACK      - server verifies that the requested address matches
  65  *                 the one it offered.
  66  *
  67  * DHCPv6 behaves pretty much the same way aside from different message names.
  68  *
  69  * Address information is embedded in either the OFFER or REQUEST message.
  70  * We chose to intercept REQUEST because this is at the last part of the
  71  * handshake and it indicates that the client intends to keep the address.
  72  * Intercepting OFFERs is unreliable because the client may receive multiple
  73  * offers from different servers, and we can't tell which address the client
  74  * will keep.
  75  *
  76  * Each DHCP message has a transaction ID. We use this transaction ID to match
  77  * REQUESTs with ACKs received from servers.
  78  *
  79  * For IPv4, the process to acquire a DHCP-assigned address is as follows:
  80  *
  81  * 1. Client sends REQUEST. a new dhcpv4_txn_t object is created and inserted
  82  *    in the the mci_v4_pending_txn table (keyed by xid). This object represents
  83  *    a new transaction. It contains the xid, the client ID and requested IP
  84  *    address.
  85  *
  86  * 2. Server responds with an ACK. The xid from this ACK is used to lookup the
  87  *    pending transaction from the mci_v4_pending_txn table. Once the object is
  88  *    found, it is removed from the pending table and inserted into the
  89  *    completed table (mci_v4_completed_txn, keyed by client ID) and the dynamic
  90  *    IP table (mci_v4_dyn_ip, keyed by IP address).
  91  *
  92  * 3. An outgoing packet that goes through the ip-nospoof path will be checked
  93  *    against the dynamic IP table. Packets that have the assigned DHCP address
  94  *    as the source IP address will pass the check and be admitted onto the
  95  *    network.
  96  *
  97  * IPv4 notes:
  98  *
  99  * If the server never responds with an ACK, there is a timer that is set after
 100  * the insertion of the transaction into the pending table. When the timer
 101  * fires, it will check whether the transaction is old (by comparing current
 102  * time and the txn's timestamp), if so the transaction will be freed. along
 103  * with this, any transaction in the completed/dyn-ip tables matching the client
 104  * ID of this stale transaction will also be freed. If the client fails to
 105  * extend a lease, we want to stop the client from using any IP addresses that
 106  * were granted previously.
 107  *
 108  * A RELEASE message from the client will not cause a transaction to be created.
 109  * The client ID in the RELEASE message will be used for finding and removing
 110  * transactions in the completed and dyn-ip tables.
 111  *
 112  *
 113  * For IPv6, the process to acquire a DHCPv6-assigned address is as follows:
 114  *
 115  * 1. Client sends REQUEST. The DUID is extracted and stored into a dhcpv6_cid_t
 116  *    structure. A new transaction structure (dhcpv6_txn_t) is also created and
 117  *    it will point to the dhcpv6_cid_t. If an existing transaction with a
 118  *    matching xid is not found, this dhcpv6_txn_t will be inserted into the
 119  *    mci_v6_pending_txn table (keyed by xid).
 120  *
 121  * 2. Server responds with a REPLY. If a pending transaction is found, the
 122  *    addresses in the reply will be placed into the dhcpv6_cid_t pointed to by
 123  *    the transaction. The dhcpv6_cid_t will then be moved to the mci_v6_cid
 124  *    table (keyed by cid). The associated addresses will be added to the
 125  *    mci_v6_dyn_ip table (while still being pointed to by the dhcpv6_cid_t).
 126  *
 127  * 3. IPv6 ip-nospoof will now check mci_v6_dyn_ip for matching packets.
```

```
 128  *      Packets with a source address matching one of the DHCPv6-assigned
 129  *      addresses will be allowed through.
 130  *
 131  * IPv6 notes:
 132  *
 133  * The v6 code shares the same timer as v4 for scrubbing stale transactions.
 134  * Just like v4, as part of removing an expired transaction, a RELEASE will be
 135  * be triggered on the cid associated with the expired transaction.
 136  *
 137  * The data structures used for v6 are slightly different because a v6 client
 138  * may have multiple addresses associated with it.
 139  */

 141 /*
 142  * These are just arbitrary limits meant for preventing abuse (e.g. a user
 143  * flooding the network with bogus transactions). They are not meant to be
 144  * user-modifiable so they are not exposed as linkprops.
 145  */
 146 static ulong_t  dhcp_max_pending_txn = 512;
 147 static ulong_t  dhcp_max_completed_txn = 512;
 148 static hrtime_t txn_cleanup_interval = 60 * NANOSEC;
  25 static time_t   txn_cleanup_interval = 60;

 150 /*
 151  * DHCPv4 transaction. It may be added to three different tables
 152  * (keyed by different fields).
 153  */
 154 typedef struct dhcpv4_txn {
 155         uint32_t                dt_xid;
 156         hrtime_t                dt_timestamp;
  33         time_t                  dt_timestamp;
 157         uint8_t                 dt_cid[DHCP_MAX_OPT_SIZE];
 158         uint8_t                 dt_cid_len;
 159         ipaddr_t                dt_ipaddr;
 160         avl_node_t              dt_node;
 161         avl_node_t              dt_ipnode;
 162         struct dhcpv4_txn       *dt_next;
 163 } dhcpv4_txn_t;
_____unchanged_portion_omitted_

 187 /*
 188  * DHCPv6 transaction. Unlike its v4 counterpart, this object gets freed up
 189  * as soon as the transaction completes or expires.
 190  */
 191 typedef struct dhcpv6_txn {
 192         uint32_t                dt_xid;
 193         hrtime_t                dt_timestamp;
  70         time_t                  dt_timestamp;
 194         dhcpv6_cid_t            *dt_cid;
 195         avl_node_t              dt_node;
 196         struct dhcpv6_txn       *dt_next;
 197 } dhcpv6_txn_t;
_____unchanged_portion_omitted_

 450 /*
 451  * Create/destroy a DHCPv4 transaction.
 452  */
 453 static dhcpv4_txn_t *
 454 create_dhcpv4_txn(uint32_t xid, uint8_t *cid, uint8_t cid_len, ipaddr_t ipaddr)
 455 {
 456         dhcpv4_txn_t    *txn;

 458         if ((txn = kmem_zalloc(sizeof (*txn), KM_NOSLEEP)) == NULL)
 459                 return (NULL);

 461         txn->dt_xid = xid;
```

```
 462         txn->dt_timestamp = gethrtime();
 339         txn->dt_timestamp = ddi_get_time();
 463         if (cid_len > 0)
 464                 bcopy(cid, &txn->dt_cid, cid_len);
 465         txn->dt_cid_len = cid_len;
 466         txn->dt_ipaddr = ipaddr;
 467         return (txn);
 468 }
_____unchanged_portion_omitted_

 505 /*
 506  * Cleanup stale DHCPv4 transactions.
 507  */
 508 static void
 509 txn_cleanup_v4(mac_client_impl_t *mcip)
 510 {
 511         dhcpv4_txn_t            *txn, *ctxn, *next, *txn_list = NULL;

 513         /*
 514          * Find stale pending transactions and place them on a list
 515          * to be removed.
 516          */
 517         for (txn = avl_first(&mcip->mci_v4_pending_txn); txn != NULL;
 518             txn = avl_walk(&mcip->mci_v4_pending_txn, txn, AVL_AFTER)) {
 519                 if (gethrtime() - txn->dt_timestamp > txn_cleanup_interval) {
 396                 if (ddi_get_time() - txn->dt_timestamp >
 397                     txn_cleanup_interval) {
 520                         DTRACE_PROBE2(found__expired__txn,
 521                             mac_client_impl_t *, mcip,
 522                             dhcpv4_txn_t *, txn);

 524                         txn->dt_next = txn_list;
 525                         txn_list = txn;
 526                 }
 527         }

 529         /*
 530          * Remove and free stale pending transactions and completed
 531          * transactions with the same client IDs as the stale transactions.
 532          */
 533         for (txn = txn_list; txn != NULL; txn = next) {
 534                 avl_remove(&mcip->mci_v4_pending_txn, txn);

 536                 ctxn = find_dhcpv4_completed_txn(mcip, txn->dt_cid,
 537                     txn->dt_cid_len);
 538                 if (ctxn != NULL) {
 539                         DTRACE_PROBE2(removing__completed__txn,
 540                             mac_client_impl_t *, mcip,
 541                             dhcpv4_txn_t *, ctxn);

 543                         remove_dhcpv4_completed_txn(mcip, ctxn);
 544                         free_dhcpv4_txn(ctxn);
 545                 }
 546                 next = txn->dt_next;
 547                 txn->dt_next = NULL;

 549                 DTRACE_PROBE2(freeing__txn, mac_client_impl_t *, mcip,
 550                     dhcpv4_txn_t *, txn);
 551                 free_dhcpv4_txn(txn);
 552         }
 553 }

 555 /*
 556  * Core logic for intercepting outbound DHCPv4 packets.
 557  */
 558 static boolean_t
```

```
 559 intercept_dhcpv4_outbound(mac_client_impl_t *mcip, ipha_t *ipha, uchar_t *end)
 560 {
 561         struct dhcp             *dh4;
 562         uchar_t                 *opt;
 563         dhcpv4_txn_t            *txn, *ctxn;
 564         ipaddr_t                ipaddr;
 565         uint8_t                 opt_len, mtype, cid[DHCP_MAX_OPT_SIZE], cid_len;
 566         mac_resource_props_t    *mrp = MCIP_RESOURCE_PROPS(mcip);

 568         if (get_dhcpv4_info(ipha, end, &dh4) != 0)
 569                 return (B_TRUE);

 571         /* ip_nospoof/allowed-ips and DHCP are mutually exclusive by default */
 572         if (allowed_ips_set(mrp, IPV4_VERSION))
 573                 return (B_FALSE);

 575         if (get_dhcpv4_option(dh4, end, CD_DHCP_TYPE, &opt, &opt_len) != 0 ||
 576             opt_len != 1) {
 577                 DTRACE_PROBE2(mtype__not__found, mac_client_impl_t *, mcip,
 578                     struct dhcp *, dh4);
 579                 return (B_TRUE);
 580         }
 581         mtype = *opt;
 582         if (mtype != REQUEST && mtype != RELEASE) {
 583                 DTRACE_PROBE3(ignored__mtype, mac_client_impl_t *, mcip,
 584                     struct dhcp *, dh4, uint8_t, mtype);
 585                 return (B_TRUE);
 586         }

 588         /* client ID is optional for IPv4 */
 589         if (get_dhcpv4_option(dh4, end, CD_CLIENT_ID, &opt, &opt_len) == 0 &&
 590             opt_len >= 2) {
 591                 bcopy(opt, cid, opt_len);
 592                 cid_len = opt_len;
 593         } else {
 594                 bzero(cid, DHCP_MAX_OPT_SIZE);
 595                 cid_len = 0;
 596         }

 598         mutex_enter(&mcip->mci_protect_lock);
 599         if (mtype == RELEASE) {
 600                 DTRACE_PROBE2(release, mac_client_impl_t *, mcip,
 601                     struct dhcp *, dh4);

 603                 /* flush any completed txn with this cid */
 604                 ctxn = find_dhcpv4_completed_txn(mcip, cid, cid_len);
 605                 if (ctxn != NULL) {
 606                         DTRACE_PROBE2(release__successful, mac_client_impl_t *,
 607                             mcip, struct dhcp *, dh4);

 609                         remove_dhcpv4_completed_txn(mcip, ctxn);
 610                         free_dhcpv4_txn(ctxn);
 611                 }
 612                 goto done;
 613         }

 615         /*
 616          * If a pending txn already exists, we'll update its timestamp so
 617          * it won't get flushed by the timer. We don't need to create new
 618          * txns for retransmissions.
 619          */
 620         if ((txn = find_dhcpv4_pending_txn(mcip, dh4->xid)) != NULL) {
 621                 DTRACE_PROBE2(update, mac_client_impl_t *, mcip,
 622                     dhcpv4_txn_t *, txn);
 623                 txn->dt_timestamp = gethrtime();
 501                 txn->dt_timestamp = ddi_get_time();
```

```
 624                 goto done;
 625         }

 627         if (get_dhcpv4_option(dh4, end, CD_REQUESTED_IP_ADDR,
 628             &opt, &opt_len) != 0 || opt_len != sizeof (ipaddr)) {
 629                 DTRACE_PROBE2(ipaddr__not__found, mac_client_impl_t *, mcip,
 630                     struct dhcp *, dh4);
 631                 goto done;
 632         }
 633         bcopy(opt, &ipaddr, sizeof (ipaddr));
 634         if ((txn = create_dhcpv4_txn(dh4->xid, cid, cid_len, ipaddr)) == NULL)
 635                 goto done;

 637         if (insert_dhcpv4_pending_txn(mcip, txn) != 0) {
 638                 DTRACE_PROBE2(insert__failed, mac_client_impl_t *, mcip,
 639                     dhcpv4_txn_t *, txn);
 640                 free_dhcpv4_txn(txn);
 641                 goto done;
 642         }
 643         start_txn_cleanup_timer(mcip);

 645         DTRACE_PROBE2(txn__pending, mac_client_impl_t *, mcip,
 646             dhcpv4_txn_t *, txn);

 648 done:
 649         mutex_exit(&mcip->mci_protect_lock);
 650         return (B_TRUE);
 651 }
_____unchanged_portion_omitted_

1112 static dhcpv6_txn_t *
1113 create_dhcpv6_txn(uint32_t xid, dhcpv6_cid_t *cid)
1114 {
1115         dhcpv6_txn_t    *txn;

1117         if ((txn = kmem_zalloc(sizeof (dhcpv6_txn_t), KM_NOSLEEP)) == NULL)
1118                 return (NULL);

1120         txn->dt_xid = xid;
1121         txn->dt_cid = cid;
1122         txn->dt_timestamp = gethrtime();
1000         txn->dt_timestamp = ddi_get_time();
1123         return (txn);
1124 }
_____unchanged_portion_omitted_

1175 /*
1176  * Cleanup stale DHCPv6 transactions.
1177  */
1178 static void
1179 txn_cleanup_v6(mac_client_impl_t *mcip)
1180 {
1181         dhcpv6_txn_t            *txn, *next, *txn_list = NULL;

1183         /*
1184          * Find stale pending transactions and place them on a list
1185          * to be removed.
1186          */
1187         for (txn = avl_first(&mcip->mci_v6_pending_txn); txn != NULL;
1188             txn = avl_walk(&mcip->mci_v6_pending_txn, txn, AVL_AFTER)) {
1189                 if (gethrtime() - txn->dt_timestamp > txn_cleanup_interval) {
1067                 if (ddi_get_time() - txn->dt_timestamp >
1068                     txn_cleanup_interval) {
1190                         DTRACE_PROBE2(found__expired__txn,
1191                             mac_client_impl_t *, mcip,
1192                             dhcpv6_txn_t *, txn);
```

```
1194                                 txn->dt_next = txn_list;
1195                                 txn_list = txn;
1196                         }
1197                 }

1199                 /*
1200                  * Remove and free stale pending transactions.
1201                  * Release any existing cids matching the stale transactions.
1202                  */
1203                 for (txn = txn_list; txn != NULL; txn = next) {
1204                         avl_remove(&mcip->mci_v6_pending_txn, txn);
1205                         release_dhcpv6_cid(mcip, txn->dt_cid);
1206                         next = txn->dt_next;
1207                         txn->dt_next = NULL;

1209                         DTRACE_PROBE2(freeing__txn, mac_client_impl_t *, mcip,
1210                             dhcpv6_txn_t *, txn);
1211                         free_dhcpv6_txn(txn);
1212                 }

1214 }

1216 /*
1217  * Core logic for intercepting outbound DHCPv6 packets.
1218  */
1219 static boolean_t
1220 intercept_dhcpv6_outbound(mac_client_impl_t *mcip, ip6_t *ip6h, uchar_t *end)
1221 {
1222         dhcpv6_message_t        *dh6;
1223         dhcpv6_txn_t            *txn;
1224         dhcpv6_cid_t            *cid = NULL;
1225         uint32_t                xid;
1226         uint8_t                 mtype;
1227         mac_resource_props_t *mrp = MCIP_RESOURCE_PROPS(mcip);

1229         if (get_dhcpv6_info(ip6h, end, &dh6) != 0)
1230                 return (B_TRUE);

1232         /* ip_nospoof/allowed-ips and DHCP are mutually exclusive by default */
1233         if (allowed_ips_set(mrp, IPV6_VERSION))
1234                 return (B_FALSE);

1236         mtype = dh6->d6m_msg_type;
1237         if (mtype != DHCPV6_MSG_REQUEST && mtype != DHCPV6_MSG_RENEW &&
1238             mtype != DHCPV6_MSG_REBIND && mtype != DHCPV6_MSG_RELEASE)
1239                 return (B_TRUE);

1241         if ((cid = create_dhcpv6_cid(dh6, end)) == NULL)
1242                 return (B_TRUE);

1244         mutex_enter(&mcip->mci_protect_lock);
1245         if (mtype == DHCPV6_MSG_RELEASE) {
1246                 release_dhcpv6_cid(mcip, cid);
1247                 goto done;
1248         }
1249         xid = DHCPV6_GET_TRANSID(dh6);
1250         if ((txn = find_dhcpv6_pending_txn(mcip, xid)) != NULL) {
1251                 DTRACE_PROBE2(update, mac_client_impl_t *, mcip,
1252                     dhcpv6_txn_t *, txn);
1253                 txn->dt_timestamp = gethrtime();
1132                 txn->dt_timestamp = ddi_get_time();
1254                 goto done;
1255         }
1256         if ((txn = create_dhcpv6_txn(xid, cid)) == NULL)
1257                 goto done;
```

```
1259         cid = NULL;
1260         if (insert_dhcpv6_pending_txn(mcip, txn) != 0) {
1261                 DTRACE_PROBE2(insert__failed, mac_client_impl_t *, mcip,
1262                     dhcpv6_txn_t *, txn);
1263                 free_dhcpv6_txn(txn);
1264                 goto done;
1265         }
1266         start_txn_cleanup_timer(mcip);

1268         DTRACE_PROBE2(txn__pending, mac_client_impl_t *, mcip,
1269             dhcpv6_txn_t *, txn);

1271 done:
1272         if (cid != NULL)
1273                 free_dhcpv6_cid(cid);

1275         mutex_exit(&mcip->mci_protect_lock);
1276         return (B_TRUE);
1277 }
```
_____**unchanged_portion_omitted_**

```
1335 /*
1336  * Timer for cleaning up stale transactions.
1337  */
1338 static void
1339 txn_cleanup_timer(void *arg)
1340 {
1341         mac_client_impl_t       *mcip = arg;

1343         mutex_enter(&mcip->mci_protect_lock);
1344         if (mcip->mci_txn_cleanup_tid == 0) {
1345                 /* do nothing if timer got cancelled */
1346                 mutex_exit(&mcip->mci_protect_lock);
1347                 return;
1348         }
1349         mcip->mci_txn_cleanup_tid = 0;

1351         txn_cleanup_v4(mcip);
1352         txn_cleanup_v6(mcip);

1354         /*
1355          * Restart timer if pending transactions still exist.
1356          */
1357         if (!avl_is_empty(&mcip->mci_v4_pending_txn) ||
1358             !avl_is_empty(&mcip->mci_v6_pending_txn)) {
1359                 DTRACE_PROBE1(restarting__timer, mac_client_impl_t *, mcip);

1361                 mcip->mci_txn_cleanup_tid = timeout(txn_cleanup_timer, mcip,
1362                     drv_usectohz(txn_cleanup_interval / (NANOSEC / MICROSEC)));
1241                     drv_usectohz(txn_cleanup_interval * 1000000));
1363         }
1364         mutex_exit(&mcip->mci_protect_lock);
1365 }

1367 static void
1368 start_txn_cleanup_timer(mac_client_impl_t *mcip)
1369 {
1370         ASSERT(MUTEX_HELD(&mcip->mci_protect_lock));
1371         if (mcip->mci_txn_cleanup_tid == 0) {
1372                 mcip->mci_txn_cleanup_tid = timeout(txn_cleanup_timer, mcip,
1373                     drv_usectohz(txn_cleanup_interval / (NANOSEC / MICROSEC)));
1252                     drv_usectohz(txn_cleanup_interval * 1000000));
1374         }
1375 }
```
_____**unchanged_portion_omitted_**