

new/usr/src/lib/libipadm/common/ipadm_addr.c

108019 Mon Jul 29 18:04:06 2013

new/usr/src/lib/libipadm/common/ipadm_addr.c

3942 inject sanity into ipadm tcp buffer size properties

3943 _snd_lowat_fraction tcp tunable has no effect

Reviewed by: Adam Leventhal <ahl@delphix.com>

Reviewed by: Peng Dai <peng.dai@delphix.com>

```
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) 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 2013 by Delphix. All rights reserved.
24 */
```

26 /*
27 * This file contains functions for address management such as creating
28 * an address, deleting an address, enabling an address, disabling an
29 * address, bringing an address down or up, setting/getting properties
30 * on an address object and listing address information
31 * for all addresses in active as well as persistent configuration.
32 */
33 #include <sys/types.h>
34 #include <sys/socket.h>
35 #include <netdb.h>
36 #include <inet/ip.h>
37 #include <string.h>
38 #include <strings.h>
39 #include <assert.h>
40 #include <sys/sockio.h>
41 #include <errno.h>
42 #include <unistd.h>
43 #include <stropts.h>
44 #include <zone.h>
45 #include <netinet/in.h>
46 #include <arpa/inet.h>
47 #include <fcntl.h>
48 #include <ctype.h>
49 #include <dhcpcagent_util.h>
50 #include <dhcpcagent_ipc.h>
51 #include <ipadm_ndpd.h>
52 #include <libdladm.h>
53 #include <libdlink.h>
54 #include <libdliptun.h>
55 #include <ifaddrs.h>
56 #include "libipadm_impl.h"
58 #define SIN6(a) ((struct sockaddr_in6 *)a)

1

new/usr/src/lib/libipadm/common/ipadm_addr.c

```
59 #define SIN(a) ((struct sockaddr_in *)a)  
61 static ipadm_status_t i_ipadm_create_addr(ipadm_handle_t, ipadm_addrobj_t,  
62                                         uint32_t);  
63 static ipadm_status_t i_ipadm_create_dhcp(ipadm_handle_t, ipadm_addrobj_t,  
64                                         uint32_t);  
65 static ipadm_status_t i_ipadm_delete_dhcp(ipadm_handle_t, ipadm_addrobj_t,  
66                                         boolean_t);  
67 static ipadm_status_t i_ipadm_get_db_addr(ipadm_handle_t, const char *,  
68                                         const char *, nvlist_t **);  
69 static ipadm_status_t i_ipadm_op_dhcp(ipadm_addrobj_t, dhcp_ipc_type_t,  
70                                         int *);  
71 static ipadm_status_t i_ipadm_validate_create_addr(ipadm_handle_t,  
72                                         ipadm_addrobj_t, uint32_t);  
73 static ipadm_status_t i_ipadm_addr_persist_nvlist(ipadm_handle_t, nvlist_t *,  
74                                         uint32_t);  
75 static ipadm_status_t i_ipadm_get_default_prefixlen(struct sockaddr_storage *,  
76                                         uint32_t *);  
77 static ipadm_status_t i_ipadm_get_static_addr_db(ipadm_handle_t,  
78                                         ipadm_addrobj_t);  
79 static boolean_t i_ipadm_is_user_aobjname_valid(const char *);  
81 /*
82  * Callback functions to retrieve property values from the kernel. These
83  * functions, when required, translate the values from the kernel to a format
84  * suitable for printing. They also retrieve DEFAULT, PERM and POSSIBLE values
85  * for a given property.
86 */
87 static ipadm_pd_getf_t i_ipadm_get_prefixlen, i_ipadm_get_addr_flag,  
88                                         i_ipadm_get_zone, i_ipadm_get_broadcast;  
90 /*
91  * Callback functions to set property values. These functions translate the
92  * values to a format suitable for kernel consumption, allocate the necessary
93  * ioctl buffers and then invoke ioctl().
94 */
95 static ipadm_pd_setf_t i_ipadm_set_prefixlen, i_ipadm_set_addr_flag,  
96                                         i_ipadm_set_zone;  
98 /* address properties description table */  
99 ipadm_prop_desc_t ipadm_addrprop_table[] = {  
100     { "broadcast", NULL, IPADMPROP_CLASS_ADDR, MOD_PROTO_NONE, 0,  
99      { "broadcast", IPADMPROP_CLASS_ADDR, MOD_PROTO_NONE, 0,  
101      NULL, NULL, i_ipadm_get_broadcast },  
103     { "deprecated", NULL, IPADMPROP_CLASS_ADDR, MOD_PROTO_NONE, 0,  
102      { "deprecated", IPADMPROP_CLASS_ADDR, MOD_PROTO_NONE, 0,  
104      i_ipadm_set_addr_flag, i_ipadm_get_onoff,  
105      i_ipadm_get_addr_flag },  
107     { "prefixlen", NULL, IPADMPROP_CLASS_ADDR, MOD_PROTO_NONE, 0,  
106      { "prefixlen", IPADMPROP_CLASS_ADDR, MOD_PROTO_NONE, 0,  
108      i_ipadm_set_prefixlen, i_ipadm_get_prefixlen,  
109      i_ipadm_get_prefixlen },  
111     { "private", NULL, IPADMPROP_CLASS_ADDR, MOD_PROTO_NONE, 0,  
110      { "private", IPADMPROP_CLASS_ADDR, MOD_PROTO_NONE, 0,  
112      i_ipadm_set_addr_flag, i_ipadm_get_onoff, i_ipadm_get_addr_flag },  
114     { "transmit", NULL, IPADMPROP_CLASS_ADDR, MOD_PROTO_NONE, 0,  
113      { "transmit", IPADMPROP_CLASS_ADDR, MOD_PROTO_NONE, 0,  
115      i_ipadm_set_addr_flag, i_ipadm_get_onoff, i_ipadm_get_addr_flag },  
117     { "zone", NULL, IPADMPROP_CLASS_ADDR, MOD_PROTO_NONE, 0,  
116      { "zone", IPADMPROP_CLASS_ADDR, MOD_PROTO_NONE, 0,  
118      i_ipadm_set_zone, NULL, i_ipadm_get_zone },
```

2

```
120     { NULL, NULL, 0, 0, 0, NULL, NULL, NULL }
119     { NULL, 0, 0, 0, NULL, NULL, NULL }
121 };

123 static ipadm_prop_desc_t up_addrprop = { "up", NULL, IPADMPROP_CLASS_ADDR,
122 static ipadm_prop_desc_t up_addrprop = { "up", IPADMPROP_CLASS_ADDR,
124                         MOD_PROTO_NONE, 0, NULL, NULL, NULL };

126 /*
127  * Helper function that initializes the 'ipadm_ifname', 'ipadm_aobjname', and
128  * 'ipadm_atype' fields of the given 'ipaddr'.
129 */
130 void
131 i_ipadm_init_addr(ipadm_addrobj_t ipaddr, const char *ifname,
132                     const char *aobjname, ipadm_addr_type_t atype)
133 {
134     bzero(ipaddr, sizeof (struct ipadm_addrobj_s));
135     (void) strlcpy(ipaddr->ipadm_ifname, ifname,
136                   sizeof (ipaddr->ipadm_ifname));
137     (void) strlcpy(ipaddr->ipadm_aobjname, aobjname,
138                   sizeof (ipaddr->ipadm_aobjname));
139     ipaddr->ipadm_atype = atype;
140 }


---

unchanged portion omitted

1374 static ipadm_prop_desc_t *
1375 i_ipadm_get_addrprop_desc(const char *pname)
1376 {
1377     int i;

1379     for (i = 0; ipadm_addrprop_table[i].ipd_name != NULL; i++) {
1380         if (strcmp(pname, ipadm_addrprop_table[i].ipd_name) == 0 ||
1381             (ipadm_addrprop_table[i].ipd_old_name != NULL &&
1382              strcmp(pname, ipadm_addrprop_table[i].ipd_old_name) == 0))
1379         if (strcmp(pname, ipadm_addrprop_table[i].ipd_name) == 0)
1383             return (&ipadm_addrprop_table[i]);
1384     }
1385     return (NULL);
1386 }


---

unchanged portion omitted
```

```
*****
54869 Mon Jul 29 18:04:08 2013
new/usr/src/lib/libipadm/common/ipadm_prop.c
3942 inject sanity into ipadm tcp buffer size properties
3943 _snd_lowat_fraction tcp tunable has no effect
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Peng Dai <peng.dai@delphix.com>
*****
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) 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 2013 by Delphix. All rights reserved.
24 */
25
26 /*
27 * This file contains routines that are used to modify/retrieve protocol or
28 * interface property values. It also holds all the supported properties for
29 * both IP interface and protocols in 'ipadm_prop_desc_t'. Following protocols
30 * are supported: IP, IPv4, IPv6, TCP, SCTP, UDP and ICMP.
31 *
32 * This file also contains walkers, which walks through the property table and
33 * calls the callback function, of the form 'ipadm_prop_wfunc_t', for every
34 * property in the table.
35 */
36
37 #include <unistd.h>
38 #include <errno.h>
39 #include <cctype.h>
40 #include <fcntl.h>
41 #include <strings.h>
42 #include <stdlib.h>
43 #include <netinet/in.h>
44 #include <arpa/inet.h>
45 #include <sys/sockio.h>
46 #include <assert.h>
47 #include <libdlink.h>
48 #include <zone.h>
49 #include "libipadm_impl.h"
50 #include <inet/tunables.h>
51
52 #define IPADM_NONESTR      "none"
53 #define DEF_METRIC_VAL     0          /* default metric value */
54
55 #define A_CNT(arr)          (sizeof (arr) / sizeof (arr[0]))
56
57 static ipadm_status_t i_ipadm_validate_if(ipadm_handle_t, const char *,
58                                         uint_t, uint_t);
```

```
60 /*
61 * Callback functions to retrieve property values from the kernel. These
62 * functions, when required, translate the values from the kernel to a format
63 * suitable for printing. For example: boolean values will be translated
64 * to on/off. They also retrieve DEFAULT, PERM and POSSIBLE values for
65 * a given property.
66 */
67 static ipadm_pd_getf_t i_ipadm_get_prop, i_ipadm_get_ifprop_flags,
68                                i_ipadm_get_mtu, i_ipadm_get_metric,
69                                i_ipadm_get_usesrc, i_ipadm_get_forwarding,
70                                i_ipadm_get_ecnsack, i_ipadm_get_hostmodel;
71
72 /*
73 * Callback function to set property values. These functions translate the
74 * values to a format suitable for kernel consumption, allocates the necessary
75 * ioctl buffers and then invokes ioctl().
76 */
77 static ipadm_pd_setf_t i_ipadm_set_prop, i_ipadm_set_mtu,
78                                i_ipadm_set_ifprop_flags,
79                                i_ipadm_set_metric, i_ipadm_set_usesrc,
80                                i_ipadm_set_forwarding, i_ipadm_set_eprivport,
81                                i_ipadm_set_ecnsack, i_ipadm_set_hostmodel;
82
83 /* array of protocols we support */
84 static int protocols[] = { MOD_PROTO_IP, MOD_PROTO_RAWIP,
85                           MOD_PROTO_TCP, MOD_PROTO_UDP,
86                           MOD_PROTO_SCTP };
87
88 /*
89 * Supported IP protocol properties.
90 */
91 static ipadm_prop_desc_t ipadm_ip_prop_table[] = {
92     { "arp", NULL, IPADMPROP_CLASS_IF, MOD_PROTO_IPV4, 0,
93       { "arp", IPADMPROP_CLASS_IF, MOD_PROTO_IPV4, 0,
94         i_ipadm_set_ifprop_flags, i_ipadm_get_onoff,
95         i_ipadm_get_ifprop_flags },
96     { "forwarding", NULL, IPADMPROP_CLASS_MODIF, MOD_PROTO_IPV4, 0,
97       { "forwarding", IPADMPROP_CLASS_MODIF, MOD_PROTO_IPV4, 0,
98         i_ipadm_set_forwarding, i_ipadm_get_onoff,
99         i_ipadm_get_forwarding },
100    { "metric", NULL, IPADMPROP_CLASS_IF, MOD_PROTO_IPV4, 0,
101      { "metric", IPADMPROP_CLASS_IF, MOD_PROTO_IPV4, 0,
102        i_ipadm_set_metric, NULL, i_ipadm_get_metric },
103    { "mtu", NULL, IPADMPROP_CLASS_IF, MOD_PROTO_IPV4, 0,
104      { "mtu", IPADMPROP_CLASS_IF, MOD_PROTO_IPV4, 0,
105        i_ipadm_set_mtu, i_ipadm_get_mtu, i_ipadm_get_mtu },
106    { "exchange_routes", NULL, IPADMPROP_CLASS_IF, MOD_PROTO_IPV4, 0,
107      { "exchange_routes", IPADMPROP_CLASS_IF, MOD_PROTO_IPV4, 0,
108        i_ipadm_set_ifprop_flags, i_ipadm_get_onoff,
109        i_ipadm_get_ifprop_flags },
110    { "usesrc", NULL, IPADMPROP_CLASS_IF, MOD_PROTO_IPV4, 0,
111      { "usesrc", IPADMPROP_CLASS_IF, MOD_PROTO_IPV4, 0,
112        i_ipadm_set_usesrc, NULL, i_ipadm_get_usesrc },
113    { "ttl", NULL, IPADMPROP_CLASS_MODULE, MOD_PROTO_IPV4, 0,
114      { "ttl", IPADMPROP_CLASS_MODULE, MOD_PROTO_IPV4, 0,
115        i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
116    { "forwarding", NULL, IPADMPROP_CLASS_MODIF, MOD_PROTO_IPV6, 0,
117      { "forwarding", IPADMPROP_CLASS_MODIF, MOD_PROTO_IPV6, 0,
```

new/usr/src/lib/libipadm/common/ipadm_prop.c

3

new/usr/src/lib/libipadm/common/ipadm_prop.c

```

170     i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
171
172     { "sack", NULL, IPADMPROP_CLASS_MODULE, MOD_PROTO_TCP, 0,
173       { "sack", IPADMPROP_CLASS_MODULE, MOD_PROTO_TCP, 0,
174         i_ipadm_set_ecnsack, i_ipadm_get_ecnsack, i_ipadm_get_ecnsack },
175
176     { "send_buf", "send_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_TCP, 0,
177       { "send_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_TCP, 0,
178         i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
179
180     { "smallest_anon_port", NULL, IPADMPROP_CLASS_MODULE, MOD_PROTO_TCP, 0,
181       { "smallest_anon_port", IPADMPROP_CLASS_MODULE, MOD_PROTO_TCP, 0,
182         i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
183
184     { "smallest_nonpriv_port", NULL, IPADMPROP_CLASS_MODULE, MOD_PROTO_TCP, 0,
185       { "smallest_nonpriv_port", IPADMPROP_CLASS_MODULE, MOD_PROTO_TCP, 0,
186         i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
187
188     { "NULL, NULL, 0, 0, 0, NULL, NULL, NULL }
189   };
190
191 /* Supported UDP protocol properties */
192 static ipadm_prop_desc_t ipadm_udp_prop_table[] = {
193   { "extra_priv_ports", NULL, IPADMPROP_CLASS_MODULE, MOD_PROTO_UDP,
194     { "extra_priv_ports", IPADMPROP_CLASS_MODULE, MOD_PROTO_UDP,
195       IPADMPROP_MULVAL, i_ipadm_set_eprivport, i_ipadm_get_prop,
196       i_ipadm_get_prop },
197
198   { "largest_anon_port", NULL, IPADMPROP_CLASS_MODULE, MOD_PROTO_UDP, 0,
199     { "largest_anon_port", IPADMPROP_CLASS_MODULE, MOD_PROTO_UDP, 0,
200       i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
201
202   { "max_buf", "max_buf", IPADMPROP_CLASS_MODULE, MOD_PROTO_UDP, 0,
203     { "recv_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_UDP, 0,
204       i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
205
206   { "recv_buf", "recv_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_UDP, 0,
207     { "send_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_UDP, 0,
208       i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
209
210   { "send_buf", "send_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_UDP, 0,
211     { "smallest_anon_port", IPADMPROP_CLASS_MODULE, MOD_PROTO_UDP, 0,
212       i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
213
214
215 /* Supported SCTP protocol properties */
216 static ipadm_prop_desc_t ipadm_sctp_prop_table[] = {
217   { "extra_priv_ports", NULL, IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP,
218     { "extra_priv_ports", IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP,
219       IPADMPROP_MULVAL, i_ipadm_set_eprivport, i_ipadm_get_prop,
220       i_ipadm_get_prop },
221
222   { "smallest_anon_port", NULL, IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP, 0,
223     { "largest_anon_port", IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP, 0,
224       i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
225
226   { "smallest_nonpriv_port", NULL, IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP, 0,
227     { "smallest_nonpriv_port", IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP, 0,
228       i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
229
230   { "NULL, NULL, 0, 0, 0, NULL, NULL, NULL }
231 }
```

```

221     i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
222     { "max_buf", "max_buf", IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP, 0,
216     { "recv_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP, 0,
224     i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
226     { "recv_buf", "recv_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP, 0,
219     { "send_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP, 0,
227     i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
228     { "send_buf", "send_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP, 0,
222     { "smallest_anon_port", IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP, 0,
230     i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
232     { "smallest_anon_port", NULL, IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP, 0,
225     { "smallest_nonpriv_port", IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP, 0,
233     i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
235     { "smallest_nonpriv_port", NULL, IPADMPROP_CLASS_MODULE, MOD_PROTO_SCTP,
236     0, i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
238     { NULL, NULL, 0, 0, 0, NULL, NULL, NULL }
228     { NULL, 0, 0, 0, NULL, NULL, NULL }
239 };

241 /* Supported ICMP protocol properties */
242 static ipadm_prop_desc_t ipadm_icmp_prop_table[] = {
243     { "max_buf", "max_buf", IPADMPROP_CLASS_MODULE, MOD_PROTO_RAWIP, 0,
233     { "recv_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_RAWIP, 0,
244     i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
246     { "recv_buf", "recv_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_RAWIP, 0,
236     { "send_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_RAWIP, 0,
247     i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
249     { "send_buf", "send_maxbuf", IPADMPROP_CLASS_MODULE, MOD_PROTO_RAWIP, 0,
250     i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop },
252     { NULL, NULL, 0, 0, 0, NULL, NULL, NULL }
239     { NULL, 0, 0, 0, NULL, NULL, NULL }
253 };

255 /*
256 * A dummy private property structure, used while handling private
257 * protocol properties (properties not yet supported by libipadm).
258 */
259 static ipadm_prop_desc_t ipadm_privprop =
260     { NULL, NULL, IPADMPROP_CLASS_MODULE, MOD_PROTO_NONE, 0,
246 static ipadm_prop_desc_t ipadm_privprop =
247     { NULL, IPADMPROP_CLASS_MODULE, MOD_PROTO_NONE, 0,
261     i_ipadm_set_prop, i_ipadm_get_prop, i_ipadm_get_prop };

263 /*
264 * Returns the property description table, for the given protocol
265 */
266 static ipadm_prop_desc_t *
267 i_ipadm_get_propdesc_table(uint_t proto)
268 {
269     switch (proto) {
270     case MOD_PROTO_IP:
271     case MOD_PROTO_IPV4:
272     case MOD_PROTO_IPV6:
273         return (ipadm_ip_prop_table);
274     case MOD_PROTO_RAWIP:
275         return (ipadm_icmp_prop_table);
276     case MOD_PROTO_TCP:

```

```

277             return (ipadm_tcp_prop_table);
278     case MOD_PROTO_UDP:
279         return (ipadm_udp_prop_table);
280     case MOD_PROTO_SCTP:
281         return (ipadm_sctp_prop_table);
282     }
284     return (NULL);
285 }

287 static ipadm_prop_desc_t *
288 i_ipadm_get_prop_desc(const char *pname, uint_t proto, int *errp)
289 {
290     int err = 0;
291     boolean_t matched_name = B_FALSE;
292     ipadm_prop_desc_t *ipdp = NULL, *ipdtbl;
294     if ((ipdtbl = i_ipadm_get_propdesc_table(proto)) == NULL) {
295         err = EINVAL;
296         goto ret;
297     }

299     for (ipdp = ipdtbl; ipdp->ipd_name != NULL; ipdp++) {
300         if (strcmp(pname, ipdp->ipd_name) == 0 ||
301             (ipdp->ipd_old_name != NULL &&
302              strcmp(pname, ipdp->ipd_old_name) == 0)) {
303             if (strcmp(pname, ipdp->ipd_name) == 0) {
304                 matched_name = B_TRUE;
305                 if (ipdp->ipd_proto == proto)
306                     break;
307             }
309             if (ipdp->ipd_name == NULL) {
310                 err = ENOENT;
311                 /* if we matched name, but failed protocol check */
312                 if (matched_name)
313                     err = EPROTO;
314                 ipdp = NULL;
315             }
316     ret:
317         if (errp != NULL)
318             *errp = err;
319         return (ipdp);
320     }

```

unchanged portion omitted

new/usr/src/lib/libipadm/common/libipadm_impl.h

8986 Mon Jul 29 18:04:10 2013

new/usr/src/lib/libipadm/common/libipadm_impl.h

3942 inject sanity into ipadm tcp buffer size properties

3943 _snd_lowat_fraction tcp tunable has no effect

Reviewed by: Adam Leventhal <ahl@delphix.com>

Reviewed by: Peng Dai <peng.dai@delphix.com>

```
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) 2010, Oracle and/or its affiliates. All rights reserved.  
23 * Copyright (c) 2013 by Delphix. All rights reserved.  
24 */  
  
26 #ifndef _LIBIPADM_IMPL_H  
27 #define _LIBIPADM_IMPL_H  
  
29 #ifdef __cplusplus  
30 extern "C" {  
31 #endif  
  
33 #include <sys/socket.h>  
34 #include <net/if.h>  
35 #include <libipadm.h>  
36 #include <libdladm.h>  
37 #include <ipadm_ipmgmt.h>  
38 #include <inet/tunables.h>  
39 #include <netinet/in.h>  
40 #include <pthread.h>  
41 #include <libinetutil.h>  
42 #include <libsocket_priv.h>  
  
44 #define IPADM_STRSIZE 256  
45 #define IPADM_ONSTR "on"  
46 #define IPADM_OFFSTR "off"  
47 #define ARP_MOD_NAME "arp"  
48 #define IPADM_LOGICAL_SEP ':'  
49 #define IPV6_MIN_MTU 1280 /* rfc2460 */  
  
51 /* mask for flags accepted by libipadm functions */  
52 #define IPADM_COMMON_OPT_MASK (IPADM_OPT_ACTIVE | IPADM_OPT_PERSIST)  
  
54 /* Opaque library handle */  
55 struct ipadm_handle {  
56     int          iph_sock;      /* socket to interface */  
57     int          iph_sock6;     /* socket to interface */  
58     int          iph_door_fd;   /* door descriptor to ipmgmtd */
```

1

new/usr/src/lib/libipadm/common/libipadm_impl.h

```
59     int          iph_rtsock;    /* routing socket */  
60     dladm_handle_t iph_dlh;    /* handle to libdladm library */  
61     uint32_t      iph_flags;    /* internal flags */  
62     pthread_mutex_t iph_lock;   /* lock to set door_fd */  
63     zoneid_t      iph_zoneid;  /* zoneid where handle was opened */  
64 };
```

unchanged portion omitted

```
94 #define ipadm_static_addr ipadm_addr_u.ipadm_static_addr_s.ipadm_addr  
95 #define ipadm_static_aname ipadm_addr_u.ipadm_static_addr_s.ipadm_ahname  
96 #define ipadm_static_prefixlen ipadm_addr_u.ipadm_static_addr_s.ipadm_prefixlen  
97 #define ipadm_static_dst_addr ipadm_addr_u.ipadm_static_addr_s.ipadm_dstaddr  
98 #define ipadm_static_dname ipadm_addr_u.ipadm_static_addr_s.ipadm_dhname  
99 #define ipadm_intfidx ipadm_addr_u.ipadm_ipv6_intfidx_s.ipadm_intfidx  
100 #define ipadm_intfidxlen ipadm_addr_u.ipadm_ipv6_intfidx_s.ipadm_intfidxlen  
101 #define ipadm_stateless ipadm_addr_u.ipadm_ipv6_intfidx_s.ipadm_stateless  
102 #define ipadm_stateful ipadm_addr_u.ipadm_ipv6_intfidx_s.ipadm_stateful  
103 #define ipadm_primary ipadm_addr_u.ipadm_dhcp_s.ipadm_primary  
104 #define ipadm_wait ipadm_addr_u.ipadm_dhcp_s.ipadm_wait  
  
106 /*  
107  * Data structures and callback functions related to property management  
108  */  
109 struct ipadm_prop_desc;  
110 typedef struct ipadm_prop_desc ipadm_prop_desc_t;  
  
112 /* property set() callback */  
113 typedef ipadm_status_t ipadm_pd_setf_t(ipadm_handle_t, const void *,  
114                                         ipadm_prop_desc_t *, const void *, uint_t, uint_t);  
  
116 /* property get() callback */  
117 typedef ipadm_status_t ipadm_pd_getf_t(ipadm_handle_t, const void *,  
118                                         ipadm_prop_desc_t *, char *, uint_t *, uint_t, uint_t);  
  
120 struct ipadm_prop_desc {  
121     char          ipd_name;      /* property name */  
122     char          ipd_old_name;  /* for backward compatibility */  
123     uint_t        ipd_class;    /* prop. class - global/perif/both */  
124     uint_t        ipd_proto;    /* protocol to which property belongs */  
125     uint_t        ipd_flags;    /* see below */  
126     ipadm_pd_setf_t *ipd_set;   /* set callback function */  
127     ipadm_pd_getf_t *ipd_get_range; /* get range callback function */  
128     ipadm_pd_getf_t *ipd_get;   /* get value callback function */  
129 };
```

unchanged portion omitted

2

new/usr/src/man/man1m/ipadm.1m

```
*****
 22256 Mon Jul 29 18:04:12 2013
new/usr/src/man/man1m/ipadm.1m
3942 inject sanity into ipadm tcp buffer size properties
3943 _snd_lowat_fraction tcp tunable has no effect
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Peng Dai <peng.dai@delphix.com>
*****
1 '\\" te
2 .\" Copyright (c) 2012, Joyent, Inc. All Rights Reserved
3 .\" Copyright (c) 2013 by Delphix. All rights reserved.
4 .\" The contents of this file are subject to the terms of the Common Development
5 .\" You can obtain a copy of the license at /usr/src/OPENSOLARIS.LICENSE or http:
6 .\" When distributing Covered Code, include this CDDL HEADER in each file and in
7 .TH IPADM 1M "May 14, 2012"
8 .SH NAME
9 ipadm \- configure IP network interfaces and protocol properties.
10 .SH SYNOPSIS
11 .LP
12 .nf
13 \fBipadm\fR create-if [\fB-t\fR] \fIinterface\fR
14 .fi
16 .LP
17 .nf
18 \fBipadm\fR disable-if [\fB-t\fR] \fIinterface\fR
19 .fi
21 .LP
22 .nf
23 \fBipadm\fR enable-if [\fB-t\fR] \fIinterface\fR
24 .fi
26 .LP
27 .nf
28 \fBipadm\fR delete-if \fIinterface\fR
29 .fi
31 .LP
32 .nf
33 \fBipadm\fR show-if [[\fB-p\fR] \fB-o\fR \fIfield\fR[,...]] [\fIinterface\fR]
34 .fi
36 .LP
37 .nf
38 \fBipadm\fR set-ifprop [\fB-t\fR] \fB-p\fR \fIprop\fR=<\fIvalue\fR[,...]> \fB-m\fR
39 .fi
41 .LP
42 .nf
43 \fBipadm\fR reset-ifprop [\fB-t\fR] \fB-p\fR \fIprop\fR \fB-m\fR \fIprotocol\fR
44 .fi
46 .LP
47 .nf
48 \fBipadm\fR show-ifprop [[\fB-c\fR]\fB-o\fR \fIfield\fR[,...]] [\fB-p\fR \fIprop
49 [\fIinterface\fR]
50 .fi
52 .LP
53 .nf
54 \fBipadm\fR create-addr [\fB-t\fR] \fB-T\fR static [\fB-d\fR]
55 \fB-a\fR {local|remote}=\fIaddr\fR[/\fIprefixlen\fR],... \fIaddrobj\fR
56 .fi
58 .LP
```

1

new/usr/src/man/man1m/ipadm.1m

```
59 .nf
60 \fBipadm\fR create-addr [\fB-t\fR] \fB-T\fR dhcp [\fB-w\fR \fIseconds\fR | forev
61 .fi
63 .LP
64 .nf
65 \fBipadm\fR create-addr [\fB-t\fR] \fB-T\fR addrconf [\fB-i\fR \fIinterface-id\f
66 [\fB-p\fR {stateful|stateless}={yes|no},... ] \fIaddrobj\fR
67 .fi
69 .LP
70 .nf
71 \fBipadm\fR down-addr [\fB-t\fR] \fIaddrobj\fR
72 .fi
74 .LP
75 .nf
76 \fBipadm\fR up-addr [\fB-t\fR] \fIaddrobj\fR
77 .fi
79 .LP
80 .nf
81 \fBipadm\fR disable-addr [\fB-t\fR] \fIaddrobj\fR
82 .fi
84 .LP
85 .nf
86 \fBipadm\fR enable-addr [\fB-t\fR] \fIaddrobj\fR
87 .fi
89 .LP
90 .nf
91 \fBipadm\fR refresh-addr [\fB-i\fR] \fIaddrobj\fR
92 .fi
94 .LP
95 .nf
96 \fBipadm\fR delete-addr [\fB-r\fR] \fIaddrobj\fR
97 .fi
99 .LP
100 .nf
101 \fBipadm\fR show-addr [[\fB-p\fR] \fB-o\fR \fIfield\fR[,...]] [\fIaddrobj\fR]
102 .fi
104 .LP
105 .nf
106 \fBipadm\fR set-addrprop [\fB-t\fR] \fB-p\fR \fIprop\fR=<\fIvalue\fR[,...]> \fIa
107 .fi
109 .LP
110 .nf
111 \fBipadm\fR reset-addrprop [\fB-t\fR] \fB-p\fR \fIprop\fR=<\fIvalue\fR[,...]> \f
112 .fi
114 .LP
115 .nf
116 \fBipadm\fR show-addrprop [[\fB-c\fR] \fB-o\fR \fIfield\fR[,...]] [\fB-p\fR \fIp
117 .fi
119 .LP
120 .nf
121 \fBipadm\fR set-prop [\fB-t\fR] \fB-p\fR \fIprop\fR[+|-]=<\fIvalue\fR[,...]> \fI
122 .fi
124 .LP
```

2

```

125 .nf
126 \fBipadm\fR reset-prop [\fB-t\fR] \fB-p\fR \fIprop\fR \fIprotocol\fR
127 .fi

129 .LP
130 .nf
131 \fBipadm\fR show-prop [\fB-c\fR] \fB-o\fR \fIfield\fR[,...] [\fB-p\fR \fIprop\
132 .fi

134 .SH DESCRIPTION
135 .sp
136 .LP

138 The \fBipadm\fR command is a stable replacement for the \fBifconfig\fR(1M) and
139 \fBndd\fR(1M) commands. It is used to create IP interfaces and to configure IP
140 addresses on those interfaces. It is also used to get, set or reset properties
141 on interfaces, addresses and protocols.
142 .LP
143 For subcommands that take an \fIaddrrobj\fR, the \fIaddrrobj\fR specifies a
144 unique address on the system. It is made up of two parts, delimited by a '/'.
145 The first part is the name of the interface and the second part is a string up
146 to 32 characters long. For example, "lo0/v4" is a loopback interface
147 addrrobj name.
148 .LP
149 For subcommands that take a \fIprotocol\fR, this can be one of
150 the following values: ip, ipv4, ipv6, icmp, tcp, sctp or udp.

152 .SH SUBCOMMANDS
153 .sp
154 .LP
155 The following subcommands are supported:
156 .sp
157 .ne 2
158 .na
159 \fB\fBcreate-if\fR [\fB-t\fR] \fIinterface\fR\fR
160 .ad
161 .sp .6
162 .RS 4n
163 The \fBcreate-if\fR subcommand is used to create an IP interface that will
164 handle both IPv4 and IPv6 packets. The interface will be enabled as part of
165 the creation process. The IPv4 interface will have the address 0.0.0.0.
166 The IPv6 interface will have the address ::.
167 .sp
168 The \fB-t\fR option (also \fB--temporary\fR) means
169 that the creation is temporary and will not be persistent across reboots.
170 .sp

172 .RE

174 .sp
175 .ne 2
176 .na
177 \fB\fBdisable-if\fR [\fB-t\fR] \fIinterface\fR\fR
178 .ad
179 .sp .6
180 .RS 4n
181 The \fBdisable-if\fR subcommand is used to disable an IP interface.
182 .sp
183 The \fB-t\fR option (also \fB--temporary\fR) means
184 that the disable is temporary and will not be persistent across reboots.
185 .sp

187 .RE

189 .sp
190 .ne 2

```

```

191 .na
192 \fB\fEnable-if\fR [\fB-t\fR] \fIinterface\fR\fR
193 .ad
194 .sp .6
195 .RS 4n
196 The \fBenable-if\fR subcommand is used to enable an IP interface.
197 .sp
198 The \fB-t\fR option (also \fB--temporary\fR) means
199 that the enable is temporary and will not be persistent across reboots.
200 .sp

202 .RE

204 .sp
205 .ne 2
206 .na
207 \fB\fBdelete-if\fR \fIinterface\fR\fR
208 .ad
209 .sp .6
210 .RS 4n
211 The \fBdelete-if\fR subcommand is used to permanently delete an IP interface.
212 .sp

214 .RE

216 .sp
217 .ne 2
218 .na
219 \fB\fBshow-if\fR [\fB-p\fR] \fB-o\fR \fIfield\fR[,...] [\fIinterface\fR\fR]
220 .ad
221 .sp .6
222 .RS 4n
223 The \fBshow-if\fR subcommand is used to show the current IP interface
224 configuration.
225 .sp
226 The \fB-p\fR option (also \fB--parsable\fR) prints
227 the output in a parsable format.
228 .sp
229 The \fB-o\fR option (also \fB--output\fR) is used
230 to select which fields will be shown. The field value can be one of the
231 following names:
232 .sp
233 .ne 2
234 .na
235 .RS 4n
236 \fBALL\fR
237 .ad
238 .RS 4n
239 Display all fields
240 .RE

242 .sp
243 .ne 2
244 .na
245 \fBIIFNAME\fR
246 .ad
247 .RS 4n
248 The name of the interface
249 .RE

251 .sp
252 .ne 2
253 .na
254 \fBSTATE\fR
255 .ad
256 .RS 4n

```

```

257 The state can be one of the following values:
258 .sp
259 .ne 2
260 .na
261 .RS 4n
262 ok - resources for the interface have been allocated
263 .sp
264 offline - the interface is offline
265 .sp
266 failed - the interface's datalink is down
267 .sp
268 down - the interface is down
269 .sp
270 disabled - the interface is disabled
271 .RE
272 .RE

274 .sp
275 .ne 2
276 .na
277 \fBCURRENT\fR
278 .ad
279 .RS 4n
280 A set of single character flags indicating the following:
281 .sp
282 .ne 2
283 .na
284 .RS 4n
285 b - broadcast (mutually exclusive with 'p')
286 .br
287 m - multicast
288 .br
289 p - point-to-point (mutually exclusive with 'b')
290 .br
291 v - virtual interface
292 .br
293 I - IPMP
294 .br
295 s - IPMP standby
296 .br
297 i - IPMP inactive
298 .br
299 V - VRRP
300 .br
301 a - VRRP accept mode
302 .br
303 4 - IPv4
304 .br
305 6 - IPv6
306 .RE
307 .RE

309 .sp
310 .ne 2
311 .na
312 \fBPERSISTENT\fR
313 .ad
314 .RS 4n
315 A set of single character flags showing what configuration will be used the
316 next time the interface is enabled:
317 .sp
318 .ne 2
319 .na
320 .RS 4n
321 s - IPMP standby
322 .br

```

```

323 4 - IPv4
324 .br
325 6 - IPv6
326 .RE
327 .RE
328 .RE

330 .RE

332 .sp
333 .ne 2
334 .na
335 \fBset-ifprop\fR [\fB-t\fR] \fB-p\fR \fIprop\fR=<\fIvalue\fR[,...]> \fB-m\fR
336 .ad
337 .sp .6
338 .RS 4n
339 The \fBset-ifprop\fR subcommand is used to set a property's value(s) on the IP
340 interface.
341 .sp
342 The \fB-t\fR option (also \fB--temporary\fR) means
343 that the setting is temporary and will not be persistent across reboots.
344 .sp
345 The \fB-p\fR option (also \fB--prop\fR) specifies the property name and
346 value(s). The property name can be one of the following:
347 .sp
348 .ne 2
349 .na

351 .RS 4n

353 \fBarp\fR
354 .ad
355 .RS 4n
356 Enables ("on") or disables ("off") ARP.
357 .RE

359 .sp
360 .ne 2
361 .na
362 \fBexchange_routes\fR
363 .ad
364 .RS 4n
365 Enables ("on") or disables ("off") the exchange of routing data.
366 .RE

368 .sp
369 .ne 2
370 .na
371 \fBforwarding\fR
372 .ad
373 .RS 4n
374 Enables ("on") or disables ("off") IP forwarding.
375 .RE

377 .sp
378 .ne 2
379 .na
380 \fBmetric\fR
381 .ad
382 .RS 4n
383 Set the routing metric to the numeric value. The value is treated as extra
384 hops to the destination.
385 .RE

387 .sp
388 .ne 2

```

```

389 .na
390 \fBmtu\fR
391 .ad
392 .RS 4n
393 Set the maximum transmission unit to the numeric value.
394 .RE

396 .sp
397 .ne 2
398 .na
399 \fBnud\fR
400 .ad
401 .RS 4n
402 Enables ("on") or disables ("off") neighbor unreachable detection.
403 .RE

405 .sp
406 .ne 2
407 .na
408 \fBusesrc\fR
409 .ad
410 .RS 4n
411 Indicates which interface to use for source address selection. A value
412 "none" may also be used.
413 .RE
414 .RE

416 .sp
417 The \fB-m\fR option (also \fB--module\fR) specifies which protocol
418 the setting applies to.
419 .sp

421 .RE
422 .RE

424 .sp
425 .ne 2
426 .na
427 \fB\fBreset-ifprop\fR [\fB-t\fR] \fB-p\fR \fIprop\fR \fB-m\fR \fIprotocol\fR \fI
428 .ad
429 .sp .6
430 .RS 4n
431 The \fBreset-ifprop\fR subcommand is used to reset an IP interface's property
432 value to the default.
433 .sp
434 The \fB-t\fR option (also \fB--temporary\fR) means
435 that the disable is temporary and will not be persistent across reboots.
436 .sp
437 The \fB-p\fR option (also \fB--prop\fR) specifies the property name.
438 See the \fBset-ifprop\fR subcommand for the list of property names.
439 .sp
440 The \fB-m\fR option (also \fB--module\fR) specifies which protocol
441 the setting applies to.
442 .sp

444 .RE

446 .sp
447 .ne 2
448 .na
449 \fB\fBshow-ifprop\fR [[\fB-c\fR]\fB-o\fR \fIfield\fR[,...]] [\fB-p\fR \fIprop\fR
450 [\fIinterface\fR]\fR
451 .ad
452 .sp .6
453 .RS 4n
454 The \fBshow-ifprop\fR subcommand is used to display the property values

```

```

455 for one or all of the IP interfaces.
456 .sp
457 The \fB-c\fR option (also \fB--parsable\fR) prints
458 the output in a parsable format.
459 .sp
460 The \fB-o\fR option (also \fB--output\fR) is used
461 to select which fields will be shown. The field value can be one of the
462 following names:
463 .sp
464 .ne 2
465 .na
466 .RS 4n
467 \fBALL\fR
468 .ad
469 .RS 4n
470 Display all fields
471 .RE

473 .sp
474 .ne 2
475 .na
476 \fBIFNAME\fR
477 .ad
478 .RS 4n
479 The name of the interface
480 .RE

482 .sp
483 .ne 2
484 .na
485 \fBPROPERTY\fR
486 .ad
487 .RS 4n
488 The name of the property
489 .RE

491 .sp
492 .ne 2
493 .na
494 \fBPROTO\fR
495 .ad
496 .RS 4n
497 The name of the protocol
498 .RE

500 .sp
501 .ne 2
502 .na
503 \fBPERM\fR
504 .ad
505 .RS 4n
506 If the property is readable ("r") and/or writable ("w").
507 .RE

509 .sp
510 .ne 2
511 .na
512 \fBCURRENT\fR
513 .ad
514 .RS 4n
515 The value of the property
516 .RE

518 .sp
519 .ne 2
520 .na

```

```

521 \fBPERISTENT\fR
522 .ad
523 .RS 4n
524 The persistent value of the property
525 .RE

527 .sp
528 .ne 2
529 .na
530 \fBDEFAULT\fR
531 .ad
532 .RS 4n
533 The default value of the property
534 .RE

536 .sp
537 .ne 2
538 .na
539 \fBPOSSIBLE\fR
540 .ad
541 .RS 4n
542 The possible values for the property
543 .RE
544 .RE

546 .sp
547 The \fB-p\fR option (also \fB--prop\fR) is used
548 to specify which properties to display. See the \fBset-ifprop\fR
549 subcommand for the list of property names.
550 .sp
551 The \fB-m\fR option (also \fB--module\fR) specifies which protocol
552 to display.
553 .sp

555 .RE

557 .sp
558 .ne 2
559 .na
560 \fB\fBcreate-addr\fR [\fB-t\fR] \fB-T\fR static [\fB-d\fR] \\
561 \fB-a\fR {local|remote}=\fIaddr\fR[/\fIprefixlen\fR],... \fIaddrobj\fR\fR
562 .br
563 \fB\fBcreate-addr\fR [\fB-t\fR] \fB-T\fR dhcp [\fB-w\fR \fIseconds\fR | forever
564 .br
565 \fB\fBcreate-addr\fR [\fB-t\fR] \fB-T\fR addrconf [\fB-i\fR \fIinterface-id\fR]
566 [\fB-p\fR {stateful|stateless}={yes|no}...] \fIaddrobj\fR\fR
567 .ad
568 .sp .6
569 .RS 4n
570 The \fBcreate-addr\fR subcommand is used to set an address on an IP interface.
571 The address will be enabled but can be disabled using the \fBdisable-addr\fR
572 subcommand. This subcommand has three different forms, depending on the
573 value of the \fB-T\fR option.
574 .sp
575 The \fB-t\fR option (also \fB--temporary\fR) means
576 that the address is temporary and will not be persistent across reboots.
577 .sp
578 The \fB-T\fR static option creates a static addrobj. This takes the following
579 options:
580 .RS 4n

582 The \fB-d\fR option (also \fB--down\fR) means the address is down.
583 .sp
584 The \fB-a\fR option (also \fB--address\fR) specifies the address.
585 The "local" or "remote" prefix can be used for a point-to-point interface.
586 In this case, both addresses must be given.

```

```

587 Otherwise, the equal sign ("=") should be omitted and the address should be
588 provided by itself and with no second address.
589 .sp
591 .RE

593 The \fB-T\fR dhcp option causes the address to be obtained via DHCP.
594 This takes the following options:
595 .RS 4n
597 The \fB-w\fR option (also \fB--wait\fR) gives the time, in seconds,
598 that the command should wait to obtain an address.
599 .sp
601 .RE

603 The \fB-T\fR addrconf option creates an auto-configured address.
604 This takes the following options:
605 .RS 4n
607 The \fB-i\fR option (also \fB--interface-id\fR) gives the interface ID to
608 be used.
609 .sp
610 The \fB-p\fR option (also \fB--prop\fR) indicates which method of
611 auto-configuration should be used.
612 .sp

614 .RE
615 .RE

617 .sp
618 .ne 2
619 .na
620 \fB\fBdown-addr\fR [\fB-t\fR] \fIaddrobj\fR\fR
621 .ad
622 .sp .6
623 .RS 4n
624 The \fBdown-addr\fR subcommand is used to stop the address. This will
625 stop packets from being sent or received.
626 .sp
627 The \fB-t\fR option (also \fB--temporary\fR) means
628 that the down is temporary and will not be persistent across reboots.
629 .sp

631 .RE

633 .sp
634 .ne 2
635 .na
636 \fB\fBup-addr\fR [\fB-t\fR] \fIaddrobj\fR\fR
637 .ad
638 .sp .6
639 .RS 4n
640 The \fBup-addr\fR subcommand is used to enable the address. This will
641 enable packets to be sent and received.
642 .sp
643 The \fB-t\fR option (also \fB--temporary\fR) means
644 that the up is temporary and will not be persistent across reboots.
645 .sp

647 .RE

649 .sp
650 .ne 2
651 .na
652 \fB\fBdisable-addr\fR [\fB-t\fR] \fIaddrobj\fR\fR

```

```

653 .ad
654 .sp .6
655 .RS 4n
656 The \fBdisable-addr\fR subcommand is used to disable the address.
657 .sp
658 The \fB-t\fR option (also \fB--temporary\fR) means
659 that the disable is temporary and will not be persistent across reboots.
660 .sp

662 .RE

664 .sp
665 .ne 2
666 .na
667 \fB\fBenable-addr\fR [\fB-t\fR] \fIaddrobj\fR\fR
668 .ad
669 .sp .6
670 .RS 4n
671 The \fBenable-addr\fR subcommand is used to enable the address.
672 .sp
673 The \fB-t\fR option (also \fB--temporary\fR) means
674 that the enable is temporary and will not be persistent across reboots.
675 .sp

677 .RE

679 .sp
680 .ne 2
681 .na
682 \fB\fBrefresh-addr\fR [\fB-i\fR] \fIaddrobj\fR\fR
683 .ad
684 .sp .6
685 .RS 4n
686 The \fBrefresh-addr\fR subcommand is used to extend the lease for DHCP
687 addresses. It also restarts duplicate address detection for Static addresses.
688 .sp
689 The \fB-i\fR option (also \fB--inform\fR) means
690 that the network configuration will be obtained from DHCP without taking
691 a lease on the address.
692 .sp

694 .RE

696 .sp
697 .ne 2
698 .na
699 \fB\fBdelete-addr\fR [\fB-r\fR] \fIaddrobj\fR\fR
700 .ad
701 .sp .6
702 .RS 4n
703 The \fBdelete-addr\fR subcommand deletes the given address.
704 .sp
705 The \fB-r\fR option (also \fB--release\fR) is used for DHCP-assigned
706 addresses to indicate that the address should be released.
707 .sp

709 .RE

711 .sp
712 .ne 2
713 .na
714 \fB\fBshow-addr\fR [(\fB-p\fR) \fB-o\fR \fIfield\fR[...]] [\fIaddrobj\fR]\fR
715 .ad
716 .sp .6
717 .RS 4n
718 The \fBshow-addr\fR subcommand is used to show the current address properties.

```

```

719 .sp
720 The \fB-p\fR option (also \fB--parsable\fR) prints
721 the output in a parsable format.
722 .sp
723 The \fB-o\fR option (also \fB--output\fR) is used
724 to select which fields will be shown. The field value can be one of the
725 following names:
726 .sp
727 .ne 2
728 .na
729 .RS 4n
730 \fBALL\fR
731 .ad
732 .RS 4n
733 Display all fields
734 .RE

736 .sp
737 .ne 2
738 .na
739 \fBADDROBJ\fR
740 .ad
741 .RS 4n
742 The name of the address
743 .RE

745 .sp
746 .ne 2
747 .na
748 \fBTYPE\fR
749 .ad
750 .RS 4n
751 The type of the address. It can be "static", "dhcp" or "addrconf".
752 .RE

754 .sp
755 .ne 2
756 .na
757 \fBSTATE\fR
758 .ad
759 .RS 4n
760 The state of the address. It can be one of the following values:
761 .sp
762 .ne 2
763 .na
764 .RS 4n
765 disabled s see the \fBdisable-addr\fR subcommand
766 .sp
767 down - see the \fBdown-addr\fR subcommand
768 .sp
769 duplicate - the address is a duplicate
770 .sp
771 inaccessible - the interface for this address has failed
772 .sp
773 ok - the address is up
774 .sp
775 tentative - duplicate address detection in progress
776 .RE
777 .RE

779 .sp
780 .ne 2
781 .na
782 \fBCURRENT\fR
783 .ad
784 .RS 4n

```

```

785 A set of single character flags indicating the following:
786 .sp
787 .ne 2
788 .na
789 .RS 4n
790 U - up
791 .br
792 u - unnumbered (matches another local address)
793 .br
794 p - private, not advertised to routing
795 .br
796 t - temporary IPv6 address
797 .br
798 d - deprecated (not used for outgoing packets)
799 .RE
800 .RE

802 .sp
803 .ne 2
804 .na
805 \fBPERSISTENT\fR
806 .ad
807 .RS 4n
808 A set of single character flags showing the configuration which will be used
809 when the address is enabled.
810 .sp
811 .ne 2
812 .na
813 .RS 4n
814 U - up
815 .br
816 p - private, not advertised to routing
817 .br
818 d - deprecated (not used for outgoing packets)
819 .RE
820 .RE

822 .sp
823 .ne 2
824 .na
825 \fBADDR\fR
826 .ad
827 .RS 4n
828 The address
829 .RE
830 .RE

832 .RE

834 .sp
835 .ne 2
836 .na
837 \fB\fBset-addrprop\fR [\fB-t\fR] \fB-p\fR \fIprop\fR=<\fIvalue\fR[,...]> \fIaddr
838 .ad
839 .sp .6
840 .RS 4n
841 The \fBset-addrprop\fR subcommand is used to set a property's value(s) on the
842 addrobj.
843 .sp
844 The \fB-t\fR option (also \fB--temporary\fR) means
845 that the setting is temporary and will not be persistent across reboots.
846 .sp
847 The \fB-p\fR option (also \fB--prop\fR) specifies the property name and
848 value(s). The property name can be one of the following:
849 .sp
850 .ne 2

```

```

851 .na
853 .RS 4n
855 \fBbroadcast\fR
856 .ad
857 .RS 4n
858 The broadcast address (read-only)
859 .RE

861 .sp
862 .ne 2
863 .na
864 \fBdeprecated\fR
865 .ad
866 .RS 4n
867 The address should not be used to send packets but can still receive packets.
868 Can be "on" or "off".
869 .RE

871 .sp
872 .ne 2
873 .na
874 \fBprefixlen\fR
875 .ad
876 .RS 4n
877 The number of bits in the IPv4 netmask or IPv6 prefix.
878 .RE

880 .sp
881 .ne 2
882 .na
883 \fBprivate\fR
884 .ad
885 .RS 4n
886 The address is not advertised to routing.
887 Can be "on" or "off".
888 .RE

890 .sp
891 .ne 2
892 .na
893 \fBtransmit\fR
894 .ad
895 .RS 4n
896 Packets can be transmitted.
897 Can be "on" or "off".
898 .RE

900 .sp
901 .ne 2
902 .na
903 \fBzone\fR
904 .ad
905 .RS 4n
906 The zone the addrobj is in.
907 .RE

909 .RE
910 .RE

912 .sp
913 .ne 2
914 .na
915 \fB\fBreset-addrprop\fR [\fB-t\fR] \fB-p\fR \fIprop\fR \fIaddrobj\fR
916 .ad

```

```

917 .sp .6
918 .RS 4n
919 The \fBreset-addrprop\fR subcommand is used to reset an addrobj's property
920 value to the default.
921 .sp
922 The \fB-t\fR option (also \fB--temporary\fR) means
923 that the disable is temporary and will not be persistent across reboots.
924 .sp
925 The \fB-p\fR option (also \fB--prop\fR) specifies the property name.
926 See the \fBset-addrprop\fR subcommand for the list of property names.
927 .sp

929 .RE

931 .sp
932 .ne 2
933 .na
934 \fB\fBshow-addrprop\fR [[\fB-c\fR]\fB-o\fR \fIfield\fR[,...]] [\fB-p\fR \fIprop\fR]
935 .ad
936 .sp .6
937 .RS 4n
938 The \fBshow-addrprop\fR subcommand is used to display the property values
939 for one or all of the addrobs.
940 .sp
941 The \fB-c\fR option (also \fB--parsable\fR) prints
942 the output in a parsable format.
943 .sp
944 The \fB-o\fR option (also \fB--output\fR) is used
945 to select which fields will be shown. The field value can be one of the
946 following names:
947 .sp
948 .ne 2
949 .na
950 .RS 4n
951 \fBALL\fR
952 .ad
953 .RS 4n
954 Display all fields
955 .RE

957 .sp
958 .ne 2
959 .na
960 \fBADDROBJ\fR
961 .ad
962 .RS 4n
963 The name of the addrobj
964 .RE

966 .sp
967 .ne 2
968 .na
969 \fBPROPERTY\fR
970 .ad
971 .RS 4n
972 The name of the property
973 .RE

975 .sp
976 .ne 2
977 .na
978 \fBPERM\fR
979 .ad
980 .RS 4n
981 If the property is readable ("r") and/or writable ("w").
982 .RE

```

```

984 .sp
985 .ne 2
986 .na
987 \fBCURRENT\fR
988 .ad
989 .RS 4n
990 The value of the property
991 .RE

993 .sp
994 .ne 2
995 .na
996 \fBPERSISTENT\fR
997 .ad
998 .RS 4n
999 The persistent value of the property
1000 .RE

1002 .sp
1003 .ne 2
1004 .na
1005 \fBDEFAULT\fR
1006 .ad
1007 .RS 4n
1008 The default value of the property
1009 .RE

1011 .sp
1012 .ne 2
1013 .na
1014 \fBPOSSIBLE\fR
1015 .ad
1016 .RS 4n
1017 The possible values for the property
1018 .RE
1019 .RE

1021 .sp
1022 The \fB-p\fR option (also \fB--prop\fR) is used
1023 to specify which properties to display. See the \fBset-addrprop\fR
1024 subcommand for the list of property names.
1025 .sp

1027 .RE

1029 .sp
1030 .ne 2
1031 .na
1032 \fB\fBset-prop\fR [\fB-t\fR] \fB-p\fR \fIprop\fR[+|-]<\fIvalue\fR[,...]> \fIpro
1033 .ad
1034 .sp .6
1035 .RS 4n
1036 The \fBset-prop\fR subcommand is used to set a property's value(s) on the
1037 protocol.
1038 .sp
1039 The \fB-t\fR option (also \fB--temporary\fR) means
1040 that the setting is temporary and will not be persistent across reboots.
1041 .sp
1042 The \fB-p\fR option (also \fB--prop\fR) specifies the property name and
1043 value(s). The optional [+|-] syntax can be used to add/remove values from the
1044 current list of values on the property.
1045 The property name can be one of the following:
1046 .sp
1047 .ne 2
1048 .na

```

```

1050 .RS 4n
1052 \fBecn\fR
1053 .ad
1054 .RS 4n
1055 Explicit congestion control (TCP-only)
1056 Can be "never", "passive" or "active".
1057 .RE

1059 \fBextra_priv_ports\fR
1060 .ad
1061 .RS 4n
1062 Additional privileged ports (SCTP, TCP or UDP)
1063 .RE

1065 \fBforwarding\fR
1066 .ad
1067 .RS 4n
1068 Packet forwarding is enabled.
1069 Can be "on" or "off".
1070 .RE

1072 \fBhoplimit\fR
1073 .ad
1074 .RS 4n
1075 The IPv6 hoplimit.
1076 .RE

1078 \fBlargest_anon_port\fR
1079 .ad
1080 .RS 4n
1081 Largest ephemeral port (SCTP, TCP or UDP)
1082 .RE

1084 \fBmax_buf\fR
1083 \fBrecv_maxbuf\fR
1085 .ad
1086 .RS 4n
1087 Maximum receive or send buffer size (ICMP, SCTP, TCP, or UDP). This also
1088 sets the upper limit for the \fBrecv_buf\fB and \fBsend_buf\fB properties.
1089 .RE

1091 \fBrecv_buf\fR
1092 .ad
1093 .RS 4n
1094 Default receive buffer size (ICMP, SCTP, TCP, or UDP). The maximum value for
1095 this property is controlled by the \fBmax_buf\fR property.
1096 .RE

1098 \fBsack\fR
1099 .ad
1100 .RS 4n
1101 Selective acknowledgement (TCP).
1102 Can be "active", "passive" or "never".
1103 .RE

1105 \fBsend_buf\fR
1096 \fBsend_maxbuf\fR
1106 .ad
1107 .RS 4n
1108 Default send buffer size (ICMP, SCTP, TCP, or UDP). The maximum value for
1109 this property is controlled by the \fBmax_buf\fR property.
1109 Send buffer size (ICMP, SCTP, TCP or UDP)
1110 .RE

```

```

1112 \fBsmallest_anon_port\fR
1113 .ad
1114 .RS 4n
1115 Smallest ephemeral port (SCTP, TCP or UDP)
1116 .RE

1118 \fBsmallest_nonpriv_port\fR
1119 .ad
1120 .RS 4n
1121 Smallest non-privileged port (SCTP, TCP or UDP)
1122 .RE

1124 \fBttl\fR
1125 .ad
1126 .RS 4n
1127 The IPv4 time-to-live.
1128 .RE

1130 .RE
1131 .RE

1133 .sp
1134 .ne 2
1135 .na
1136 \fB\fBreset-prop\fR [\fB-t\fR] \fB-p\fR \fIprop\fR \fIprotocol\fR\fR
1137 .ad
1138 .sp .6
1139 .RS 4n
1140 The \fBreset-prop\fR subcommand is used to reset a protocol's property
1141 value to the default.
1142 .sp
1143 The \fB-t\fR option (also \fB--temporary\fR) means
1144 that the disable is temporary and will not be persistent across reboots.
1145 .sp
1146 The \fB-p\fR option (also \fB--prop\fR) specifies the property name.
1147 See the \fBset-prop\fR subcommand for the list of property names.
1148 .sp

1150 .RE

1152 .sp
1153 .ne 2
1154 .na
1155 \fB\fBshow-prop\fR [[\fB-c\fR]\fB-o\fR \fIfield\fR[,...]] [\fB-p\fR \fIprop\fR,..
1156 .ad
1157 .sp .6
1158 .RS 4n
1159 The \fBshow-prop\fR subcommand is used to display the property values
1160 for one or all of the protocols.
1161 .sp
1162 The \fB-c\fR option (also \fB--parsable\fR) prints
1163 the output in a parsable format.
1164 .sp
1165 The \fB-o\fR option (also \fB--output\fR) is used
1166 to select which fields will be shown. The field value can be one of the
1167 following names:
1168 .sp
1169 .ne 2
1170 .na
1171 .RS 4n
1172 \fBALL\fR
1173 .ad
1174 .RS 4n
1175 Display all fields
1176 .RE

```

```

1178 .sp
1179 .ne 2
1180 .na
1181 \fBPROTO\fR
1182 .ad
1183 .RS 4n
1184 The name of the protocol
1185 .RE

1187 .sp
1188 .ne 2
1189 .na
1190 \fBPROPERTY\fR
1191 .ad
1192 .RS 4n
1193 The name of the property
1194 .RE

1196 .sp
1197 .ne 2
1198 .na
1199 \fBPERM\fR
1200 .ad
1201 .RS 4n
1202 If the property is readable ("r") and/or writable ("w").
1203 .RE

1205 .sp
1206 .ne 2
1207 .na
1208 \fBCURRENT\fR
1209 .ad
1210 .RS 4n
1211 The value of the property
1212 .RE

1214 .sp
1215 .ne 2
1216 .na
1217 \fBPERSISTENT\fR
1218 .ad
1219 .RS 4n
1220 The persistent value of the property
1221 .RE

1223 .sp
1224 .ne 2
1225 .na
1226 \fBDEFAULT\fR
1227 .ad
1228 .RS 4n
1229 The default value of the property
1230 .RE

1232 .sp
1233 .ne 2
1234 .na
1235 \fBPOSSIBLE\fR
1236 .ad
1237 .RS 4n
1238 The possible values for the property
1239 .RE
1240 .RE

1242 .sp

```

```

1243 The \fB-p\fR option (also \fB--prop\fR) is used
1244 to specify which properties to display. See the \fBset-prop\fR
1245 subcommand for the list of property names.
1246 .sp
1248 .RE
1250 .SH SEE ALSO
1251 .sp
1252 .LP
1253 \fBifconfig\fR(1M), \fBdladm\fR(1M), \fBndd\fR(1M), \fBzonecfg\fR(1M),
1254 \fBarp\fR(1M), \fBcfgadm\fR(1M), \fBif_mpadm\fR(1M), \fBnsswitch.conf\fR(4),
1255 and \fBdhcp\fR(5).

```

new/usr/src/uts/common/inet/ip/icmp.c

153451 Mon Jul 29 18:04:14 2013

new/usr/src/uts/common/inet/ip/icmp.c

3942 inject sanity into ipadmin tcp buffer size properties

3943 _snd_lowat_fraction tcp tunable has no effect

Reviewed by: Adam Leventhal <ahl@delphix.com>

Reviewed by: Peng Dai <peng.dai@delphix.com>

```
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) 1991, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 2013 by Delphix. All rights reserved.
24 */
25 /* Copyright (c) 1990 Mentor Inc. */
```

```
27 #include <sys/types.h>
28 #include <sys/stream.h>
29 #include <sys/stropts.h>
30 #include <sys/strlog.h>
31 #include <sys/strsun.h>
32 #define _SUN_TPI_VERSION 2
33 #include <sys/tihdr.h>
34 #include <sys/timod.h>
35 #include <sys/ddi.h>
36 #include <sys/sunddi.h>
37 #include <sys/strsnbr.h>
38 #include <sys/suntpi.h>
39 #include <sys/xti_inet.h>
40 #include <sys/cmn_err.h>
41 #include <sys/kmem.h>
42 #include <sys/cred.h>
43 #include <sys/policy.h>
44 #include <sys/priv.h>
45 #include <sys/ucred.h>
46 #include <sys/zone.h>

48 #include <sys/sockio.h>
49 #include <sys/socket.h>
50 #include <sys/socketvar.h>
51 #include <sys/vtrace.h>
52 #include <sys/sdt.h>
53 #include <sys/debug.h>
54 #include <sys/isa_defs.h>
55 #include <sys/random.h>
56 #include <netinet/in.h>
57 #include <netinet/ip6.h>
58 #include <netinet/icmp6.h>
```

1

new/usr/src/uts/common/inet/ip/icmp.c

```
59 #include <netinet/udp.h>
60
61 #include <inet/common.h>
62 #include <inet/ip.h>
63 #include <inet/ip_impl.h>
64 #include <inet/ipsec_impl.h>
65 #include <inet/ip6.h>
66 #include <inet/ip_ire.h>
67 #include <inet/ip_if.h>
68 #include <inet/ip_multi.h>
69 #include <inet/ip_ndp.h>
70 #include <inet/proto_set.h>
71 #include <inet/mib2.h>
72 #include <inet/nd.h>
73 #include <inet/optcom.h>
74 #include <inet/snmpcom.h>
75 #include <inet/kstatcom.h>
76 #include <inet/ipclassifier.h>
77
78 #include <sys/tsol/label.h>
79 #include <sys/tsol/tnet.h>
80
81 #include <inet/rawip_impl.h>
82
83 #include <sys/disp.h>
84
85 /*
86  * Synchronization notes:
87  *
88  * RAWIP is MT and uses the usual kernel synchronization primitives. We use
89  * conn_lock to protect the icmp_t.
90  *
91  * Plumbing notes:
92  * ICMP is always a device driver. For compatibility with mibopen() code
93  * it is possible to I_PUSH "icmp", but that results in pushing a passthrough
94  * dummy module.
95 */
96 static void icmp_addr_req(queue_t *q, mblk_t *mp);
97 static void icmp_tpi_bind(queue_t *q, mblk_t *mp);
98 static void icmp_bind_proto(icmp_t *icmp);
99 static int icmp_build_hdr_template(conn_t *, const in6_addr_t *,
100     const in6_addr_t *, uint32_t);
101 static void icmp_capability_req(queue_t *q, mblk_t *mp);
102 static int icmp_close(queue_t *q, int flags);
103 static void icmp_close_free(conn_t *);
104 static void icmp_tpi_connect(queue_t *q, mblk_t *mp);
105 static void icmp_tpi_disconnect(queue_t *q, mblk_t *mp);
106 static void icmp_err_ack(queue_t *q, mblk_t *mp, t_scalar_t t_error,
107     int sys_error);
108 static void icmp_err_ack_prim(queue_t *q, mblk_t *mp, t_scalar_t primitive,
109     t_scalar_t tl ierr, int sys_error);
110 static void icmp_icmp_input(void *arg1, mblk_t *mp, void *arg2,
111     ip_recv_attr_t *);
112 static void icmp_icmp_error_ipv6(conn_t *connp, mblk_t *mp,
113     ip_recv_attr_t *);
114 static void icmp_info_req(queue_t *q, mblk_t *mp);
115 static void icmp_input(void *, mblk_t *, void *, ip_recv_attr_t *);
116 static conn_t *icmp_open(int family, cred_t *credp, int *err, int flags);
117 static int icmp_openv4(queue_t *q, dev_t *devp, int flag, int sflag,
118     cred_t *credp);
119 static int icmp_openv6(queue_t *q, dev_t *devp, int flag, int sflag,
120     cred_t *credp);
121 static boolean_t icmp_opt_allow_udr_set(t_scalar_t level, t_scalar_t name);
122 int icmp_opt_set(conn_t *connp, uint_t optset_context,
123     int level, int name, uint_t inlen,
```

2

```

124         uchar_t *invalp, uint_t *outlenp, uchar_t *outvalp,
125         void *thisdgAttrs, cred_t *cr);
126 int         icmp_opt_get(conn_t *connp, int level, int name,
127                         uchar_t *ptr);
128 static int   icmp_output_newdst(conn_t *connp, mblk_t *data_mp, sin_t *sin,
129                                 sin6_t *sin6, cred_t *cr, pid_t pid, ip_xmit_attr_t *ixa);
130 static mblk_t *icmp_prepend_hdr(conn_t *, ip_xmit_attr_t *, const ip_pkt_t *,
131                                 const in6_addr_t *, const in6_addr_t *, uint32_t, mblk_t *, int *);
132 static mblk_t *icmp_prepend_header_template(conn_t *, ip_xmit_attr_t *,
133                                 mblk_t *, const in6_addr_t *, uint32_t, int *);
134 static int   icmp_snmp_set(queue_t *q, t_scalar_t level, t_scalar_t name,
135                         uchar_t *ptr, int len);
136 static void  icmp_ud_err(queue_t *q, mblk_t *mp, t_scalar_t err);
137 static void  icmp_tp1_unbind(queue_t *q, mblk_t *mp);
138 static void  icmp_wput(queue_t *q, mblk_t *mp);
139 static void  icmp_wput_fallback(queue_t *q, mblk_t *mp);
140 static void  icmp_wput_other(queue_t *q, mblk_t *mp);
141 static void  icmp_wput_iodata(queue_t *q, mblk_t *mp);
142 static void  icmp_wput_restricted(queue_t *q, mblk_t *mp);
143 static void  icmp_upl_recv(conn_t *, mblk_t *, uint_t);

145 static void  *rawip_stack_init(netstackid_t stackid, netstack_t *ns);
146 static void  rawip_stack_fini(netstackid_t stackid, void *arg);

148 static void  *rawip_kstat_init(netstackid_t stackid);
149 static void  rawip_kstat_fini(netstackid_t stackid, kstat_t *ksp);
150 static int   rawip_kstat_update(kstat_t *kp, int rw);
151 static void  rawip_stack_shutdown(netstackid_t stackid, void *arg);

153 /* Common routines for TPI and socket module */
154 static conn_t *rawip_do_open(int, cred_t *, int *, int);
155 static void  rawip_do_close(conn_t *);
156 static int   rawip_do_bind(conn_t *, struct sockaddr *, socklen_t);
157 static int   rawip_do_unbind(conn_t *);
158 static int   rawip_do_connect(conn_t *, const struct sockaddr *, socklen_t,
159                           cred_t *, pid_t);

161 int         rawip_getsockname(sock_lower_handle_t, struct sockaddr *,
162                             socklen_t *, cred_t *);
163 int         rawip_getpeername(sock_lower_handle_t, struct sockaddr *,
164                             socklen_t *, cred_t *);

166 static struct module_info icmp_mod_info = {
167     5707, "icmp", 1, INFPSZ, 512, 128
168 };


---


unchanged_portion_omitted

216 static int   icmp_set_buf_prop(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
217                               const char *ifname, const void *pval, uint_t flags)
218 {
219     return (mod_set_buf_prop(stack->netstack_icmp->is_propinfo_tbl,
220                           stack, cr, pinfo, ifname, pval, flags));
221 }

224 static int   icmp_get_buf_prop(netstack_t *stack, mod_prop_info_t *pinfo, const char *ifname,
225                               void *val, uint_t psize, uint_t flags)
226 {
227     return (mod_get_buf_prop(stack->netstack_icmp->is_propinfo_tbl, stack,
228                           pinfo, ifname, val, psize, flags));
229 }

232 /*
233 * All of these are alterable, within the min/max values given, at run time.
234 */

```

```

235     * Note: All those tunables which do not start with "icmp_" are Committed and
236     * therefore are public. See PSARC 2010/080.
237     */
238 static mod_prop_info_t icmp_propinfo_tbl[] = {
239     /* tunable - 0 */
240     { "_wroff_extra", MOD_PROTO_RAWIP,
241       mod_set_uint32, mod_get_uint32,
242       {0, 128, 32}, {32} },
243
244     { "_ipv4_ttl", MOD_PROTO_RAWIP,
245       mod_set_uint32, mod_get_uint32,
246       {1, 255, 255}, {255} },
247
248     { "_ipv6_hoplimit", MOD_PROTO_RAWIP,
249       mod_set_uint32, mod_get_uint32,
250       {0, IPV6_MAX_HOPS, IPV6_DEFAULT_HOPS},
251       {IPV6_DEFAULT_HOPS} },
252
253     { "_bsd_compat", MOD_PROTO_RAWIP,
254       mod_set_boolean, mod_get_boolean,
255       {B_TRUE}, {B_TRUE} },
256
257     { "send_buf", MOD_PROTO_RAWIP,
258       icmp_set_buf_prop, icmp_get_buf_prop,
259     { "send_maxbuf", MOD_PROTO_RAWIP,
260       mod_set_uint32, mod_get_uint32,
261       {4096, 65536, 8192}, {8192} },
262
263     { "_xmit_lowat", MOD_PROTO_RAWIP,
264       mod_set_uint32, mod_get_uint32,
265       {0, 65536, 1024}, {1024} },
266
267     { "recv_buf", MOD_PROTO_RAWIP,
268       icmp_set_buf_prop, icmp_get_buf_prop,
269     { "recv_maxbuf", MOD_PROTO_RAWIP,
270       mod_set_uint32, mod_get_uint32,
271       {4096, 65536, 8192}, {8192} },
272
273     { "max_buf", MOD_PROTO_RAWIP,
274     { "max_buf", MOD_PROTO_RAWIP,
275       mod_set_uint32, mod_get_uint32,
276       {65536, ULP_MAX_BUF, 256*1024}, {256*1024} },
277       {65536, 1024*1024, 256*1024}, {256 * 1024} },
278
279     { "_pmtu_discovery", MOD_PROTO_RAWIP,
280       mod_set_boolean, mod_get_boolean,
281       {B_FALSE}, {B_FALSE} },
282
283     { "_sendto_ignerr", MOD_PROTO_RAWIP,
284       mod_set_boolean, mod_get_boolean,
285       {B_FALSE}, {B_FALSE} },
286
287     { "?", MOD_PROTO_RAWIP, NULL, mod_get_allprop, {0}, {0} },
288     { NULL, 0, NULL, NULL, {0}, {0} }
289 };


---


unchanged_portion_omitted

```

new/usr/src/uts/common/inet/ip/ip_if.c

533173 Mon Jul 29 18:04:15 2013

new/usr/src/uts/common/inet/ip/ip_if.c

3942 inject sanity into ipadv tcp buffer size properties

3943 _snd_lowat_fraction tcp tunable has no effect

Reviewed by: Adam Leventhal <ahl@delphix.com>

Reviewed by: Peng Dai <peng.dai@delphix.com>

```
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) 1991, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 1990 Mentor Inc.
24 * Copyright (c) 2013 by Delphix. All rights reserved.
25 */
```

```
27 /*
28 * This file contains the interface control functions for IP.
29 */
```

```
31 #include <sys/types.h>
32 #include <sys/stream.h>
33 #include <sys/dlpi.h>
34 #include <sys/stropts.h>
35 #include <sys/strsun.h>
36 #include <sys/sysmacros.h>
37 #include <sys/strsubr.h>
38 #include <sys/strlog.h>
39 #include <sys/ddi.h>
40 #include <sys/sunddi.h>
41 #include <sys/cmn_err.h>
42 #include <sys/kstat.h>
43 #include <sys/debug.h>
44 #include <sys/zone.h>
45 #include <sys/sunldi.h>
46 #include <sys/file.h>
47 #include <sys/bitmap.h>
48 #include <sys/cpuvar.h>
49 #include <sys/time.h>
50 #include <sys/ctype.h>
51 #include <sys/kmem.h>
52 #include <sys/sysdm.h>
53 #include <sys/param.h>
54 #include <sys/socket.h>
55 #include <sys/isa_defs.h>
56 #include <net/if.h>
57 #include <net/if_arp.h>
58 #include <net/if_types.h>
```

1

new/usr/src/uts/common/inet/ip/ip_if.c

```
59 #include <net/if_dl.h>
60 #include <net/route.h>
61 #include <sys/sockio.h>
62 #include <netinet/in.h>
63 #include <netinet/ip6.h>
64 #include <netinet/icmp6.h>
65 #include <netinet/igmp_var.h>
66 #include <sys/policy.h>
67 #include <sys/ethernet.h>
68 #include <sys/callb.h>
69 #include <sys/md5.h>

71 #include <inet/common.h> /* for various inet/mi.h and inet/nd.h needs */
72 #include <inet/mi.h>
73 #include <inet/nd.h>
74 #include <inet/tunables.h>
75 #include <inet/arp.h>
76 #include <inet/ip_arp.h>
77 #include <inet/mib2.h>
78 #include <inet/ip.h>
79 #include <inet/ip6.h>
80 #include <inet/ip6_asp.h>
81 #include <inet/tcp.h>
82 #include <inet/ip_multi.h>
83 #include <inet/ip_ire.h>
84 #include <inet/ip_ftable.h>
85 #include <inet/ip_rts.h>
86 #include <inet/ip_ndp.h>
87 #include <inet/ip_if.h>
88 #include <inet/ip_impl.h>
89 #include <inet/sctp_ip.h>
90 #include <inet/ip_netinfo.h>
91 #include <inet/ilb_ip.h>

93 #include <netinet/igmp.h>
94 #include <inet/ip_listutils.h>
95 #include <inet/ipclassifier.h>
96 #include <sys/mac_client.h>
97 #include <sys/dld.h>
98 #include <sys/mac_flow.h>

100 #include <sys/systeminfo.h>
101 #include <sys/bootconf.h>

103 #include <sys/tsol/tndb.h>
104 #include <sys/tsol/tnt.h>

106 #include <inet/rawip_impl.h> /* needed for icmp_stack_t */
107 #include <inet/udp_impl.h> /* needed for udp_stack_t */

109 /* The character which tells where the ill_name ends */
110 #define IPIF_SEPARATOR_CHAR ':'

112 /* IP ioctl function table entry */
113 typedef struct ipft_s {
114     int      ipft_cmd;
115     pfi_t   ipft_pfi;
116     int      ipft_min_size;
117     int      ipft_flags;
118 } ipft_t;
119 _____  
unchanged_portion_omitted_____  
119 _____
```

```
8846 /*
8847 * process the SIOC{SET|GET}PROP ioctl's
8848 */
8849 /* ARGSUSED */
```

2

```

8850 static void
8851 ip_ioctl_getsetprop(queue_t *q, mblk_t *mp)
8852 {
8853     struct iocblk    *iocp = (struct iocblk *)mp->b_rptr;
8854     mblk_t          *mpl = mp->b_cont;
8855     mod_ioc_prop_t *pioc;
8856     mod_prop_info_t *ptbl = NULL, *pinfo = NULL;
8857     ip_stack_t      *ipst;
8858     icmp_stack_t    *is;
8859     tcp_stack_t     *tcps;
8860     sctp_stack_t    *sctps;
8861     udp_stack_t     *us;
8862     netstack_t      *stack;
8863     void            *cbarg;
8864     cred_t          *cr;
8865     boolean_t        set;
8866     int              err;
8867
8868     ASSERT(q->q_next == NULL);
8869     ASSERT(CONN_Q(q));
8870
8871     if (!getset_ioctl_checks(mp)) {
8872         miocnak(q, mp, 0, EINVAL);
8873         return;
8874     }
8875     ipst = CONNQ_TO_IPST(q);
8876     stack = ipst->ips_netstack;
8877     pioc = (mod_ioc_prop_t *)mpl->b_rptr;
8878
8879     switch (pioc->mpr_proto) {
8880     case MOD_PROTO_IP:
8881     case MOD_PROTO_IPV4:
8882     case MOD_PROTO_IPV6:
8883         ptbl = ipst->ips_propinfo_tbl;
8884         cbarg = ipst;
8885         break;
8886     case MOD_PROTO_RAWIP:
8887         ptbl = stack->netstack_icmp->is_propinfo_tbl;
8888         is = stack->netstack_icmp;
8889         ptbl = is->is_propinfo_tbl;
8890         cbarg = is;
8891         break;
8892     case MOD_PROTO_TCP:
8893         ptbl = stack->netstack_tcp->tcps_propinfo_tbl;
8894         tcps = stack->netstack_tcp;
8895         ptbl = tcps->tcps_propinfo_tbl;
8896         cbarg = tcps;
8897         break;
8898     case MOD_PROTO_UDP:
8899         ptbl = stack->netstack_udp->us_propinfo_tbl;
8900         us = stack->netstack_udp;
8901         ptbl = us->us_propinfo_tbl;
8902         cbarg = us;
8903         break;
8904     case MOD_PROTO_SCTP:
8905         ptbl = stack->netstack_sctp->sctps_propinfo_tbl;
8906         sctps = stack->netstack_sctp;
8907         ptbl = sctps->sctps_propinfo_tbl;
8908         cbarg = sctps;
8909         break;
8910     default:
8911         miocnak(q, mp, 0, EINVAL);
8912         return;
8913     }
8914
8915     pinfo = mod_prop_lookup(ptbl, pioc->mpr_name, pioc->mpr_proto);

```

```

8916     if (pinfo == NULL) {
8917         /* search for given property in respective protocol propinfo table */
8918         for (pinfo = ptbl; pinfo->mpi_name != NULL; pinfo++) {
8919             if (strcmp(pinfo->mpi_name, pioc->mpr_name) == 0 &&
8920                 pinfo->mpr_proto == pioc->mpr_proto)
8921                 break;
8922         }
8923         if (pinfo->mpi_name == NULL) {
8924             miocnak(q, mp, 0, ENOENT);
8925             return;
8926         }
8927
8928         set = (iocp->ioc_cmd == SIOCSETPROP) ? B_TRUE : B_FALSE;
8929         if (set && pinfo->mpi_setf != NULL) {
8930             cr = msg_getcred(mp, NULL);
8931             if (cr == NULL)
8932                 cr = iocp->ioc_cr;
8933             err = pinfo->mpi_setf(stack, cr, pinfo, pioc->mpr_ifname,
8934                                     err = pinfo->mpi_setf(cbarg, cr, pinfo, pioc->mpr_ifname,
8935                                     pioc->mpr_val, pioc->mpr_flags);
8936         } else if (!set && pinfo->mpi_getf != NULL) {
8937             err = pinfo->mpi_getf(stack, pinfo, pioc->mpr_ifname,
8938             err = pinfo->mpi_getf(cbarg, pinfo, pioc->mpr_ifname,
8939             pioc->mpr_val, pioc->mpr_valsize, pioc->mpr_flags);
8940         } else {
8941             err = EPERM;
8942         }
8943
8944         if (err != 0) {
8945             miocnak(q, mp, 0, err);
8946         } else {
8947             if (set)
8948                 miocack(q, mp, 0, 0);
8949             else /* For get, we need to return back the data */
8950                 miocack(q, mp, iocp->ioc_count, 0);
8951         }
8952
8953     /*
8954      * process the legacy ND_GET, ND_SET ioctl just for {ip|ip6}_forwarding
8955      * as several routing daemons have unfortunately used this 'unpublished'
8956      * but well-known ioctls.
8957      */
8958     /* ARGUSED */
8959     static void
8960     ip_process_legacy_nddprop(queue_t *q, mblk_t *mp)
8961     {
8962         struct iocblk    *iocp = (struct iocblk *)mp->b_rptr;
8963         mblk_t          *mpl = mp->b_cont;
8964         char            *pname, *pval, *buf;
8965         uint_t          bufsize, proto;
8966         mod_prop_info_t *pinfo = NULL;
8967         mod_prop_info_t *ptbl = NULL, *pinfo = NULL;
8968         ip_stack_t      *ipst;
8969         int              err = 0;
8970
8971         ASSERT(CONN_Q(q));
8972         ipst = CONNQ_TO_IPST(q);
8973
8974         if (iocp->ioc_count == 0 || mpl == NULL) {
8975             miocnak(q, mp, 0, EINVAL);
8976             return;
8977         }
8978
8979         mpl->b_datap->db_lim[-1] = '\0'; /* Force null termination */
8980         pval = buf = pname = (char *)mpl->b_rptr;

```

```
8954     bufsize = MBLKL(mp1);
8956     if (strcmp(pname, "ip_forwarding") == 0) {
8957         pname = "forwarding";
8958         proto = MOD_PROTO_IPV4;
8959     } else if (strcmp(pname, "ip6_forwarding") == 0) {
8960         pname = "forwarding";
8961         proto = MOD_PROTO_IPV6;
8962     } else {
8963         miocnak(q, mp, 0, EINVAL);
8964         return;
8965     }
8966     pinfo = mod_prop_lookup(ipst->ips_propinfo_tbl, pname, proto);
8967     ptbl = ipst->ips_propinfo_tbl;
8968     for (pinfo = ptbl; pinfo->mpi_name != NULL; pinfo++) {
8969         if (strcmp(pinfo->mpi_name, pname) == 0 &&
8970             pinfo->mpi_proto == proto)
8971             break;
8972     }
8973     ASSERT(pinfo->mpi_name != NULL);
8974     switch (iocp->ioc_cmd) {
8975     case ND_GET:
8976         if ((err = pinfo->mpi_getf(ipst->ips_netstack, pinfo, NULL, buf,
8977                                     bufsize, 0)) == 0) {
8978             if ((err = pinfo->mpi_getf(ipst, pinfo, NULL, buf, bufsize,
8979                                         0)) == 0) {
8980                 miocack(q, mp, iocp->ioc_count, 0);
8981                 return;
8982             }
8983             break;
8984     case ND_SET:
8985         /*
8986         * buffer will have property name and value in the following
8987         * format,
8988         * <property name>'\'0'<property value>'\'0', extract them;
8989         */
8990         while (*pval++) {
8991             noop;
8992             if (!*pval || pval >= (char *)mp1->b_wptr) {
8993                 err = EINVAL;
8994             } else if ((err = pinfo->mpi_setf(ipst->ips_netstack, NULL,
8995                                             pinfo, NULL, pval, 0)) == 0) {
8996             } else if ((err = pinfo->mpi_setf(ipst, NULL, pinfo, NULL,
8997                                             pval, 0)) == 0) {
8998                 miocack(q, mp, 0, 0);
8999             }
9000         }
9001         break;
9002     default:
9003         err = EINVAL;
9004         break;
9005     }
9006     miocnak(q, mp, 0, err);
9007 }
```

unchanged_portion_omitted_

```
*****
27327 Mon Jul 29 18:04:18 2013
new/usr/src/uts/common/inet/ip/ip_tunables.c
3942 inject sanity into ipadmin tcp buffer size properties
3943 _snd_lowat_fraction tcp tunable has no effect
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Peng Dai <peng.dai@delphix.com>
*****
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) 1991, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 2013 by Delphix. All rights reserved.
24 */
25 /* Copyright (c) 1990 Mentor Inc. */

27 #include <inet/ip.h>
28 #include <inet/ip6.h>
29 #include <inet/ip_if.h>
30 #include <inet/ip_ire.h>
31 #include <inet/ipclassifier.h>
32 #include <inet/ip_impl.h>
33 #include <inet/tunables.h>
34 #include <sys/sunddi.h>
35 #include <sys/policy.h>

37 /* How long, in seconds, we allow frags to hang around. */
38 #define IP_REASM_TIMEOUT 15
39 #define IPV6_REASM_TIMEOUT 60

41 /*
42 * Set ip{,6}_forwarding values. If the value is being set on an ill,
43 * find the ill and set the value on it. On the other hand if we are modifying
44 * global property, modify the global value and set the value on all the ills.
45 */
46 /* ARGSUSED */
47 static int
48 ip_set_forwarding(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
49 ip_set_forwarding(void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
50 const char *ifname, const void* pval, uint_t flags)
51 {
52     char *end;
53     unsigned long new_value;
54     boolean_t per_ill, isv6;
55     ill_walk_context_t *ill;
56     ip_stack_t *ipst = stack->netstack_ip;
57     ip_stack_t *ipst = (ip_stack_t *)cbarg;
```

```
58     if (flags & MOD_PROP_DEFAULT) {
59         new_value = pinfo->prop_def_bval;
60     } else {
61         if (ddi strtoul(pval, &end, 10, &new_value) != 0 ||
62             *end != '\0')
63             return (EINVAL);
64         if (new_value != B_TRUE && new_value != B_FALSE)
65             return (EINVAL);
66     }
67
68     per_ill = (ifname != NULL && ifname[0] != '\0');
69     /*
70      * if it's not per ill then set the global property and bring all the
71      * ills up to date with the new global value.
72      */
73     if (!per_ill)
74         pinfo->prop_cur_bval = (new_value == 1 ? B_TRUE : B_FALSE);
75
76     isv6 = (pinfo->mpi_proto == MOD_PROTO_IPV6 ? B_TRUE : B_FALSE);
77     rw_enter(&ipst->ips_ill_g_lock, RW_READER);
78     if (isv6)
79         ill = ILL_START_WALK_V6(&ctx, ipst);
80     else
81         ill = ILL_START_WALK_V4(&ctx, ipst);
82
83     for (; ill != NULL; ill = ill_next(&ctx, ill)) {
84         /*
85          * if the property needs to be set on a particular
86          * interface, look for that interface.
87          */
88         if (per_ill && strcmp(ifname, ill->ill_name) != 0)
89             continue;
90         (void) ill_forward_set(ill, new_value != 0);
91     }
92     rw_exit(&ipst->ips_ill_g_lock);
93
94     return (0);
95 }

96 static int
97 ip_get_forwarding(netstack_t *stack, mod_prop_info_t *pinfo, const char *ifname,
98 ip_get_forwarding(void *cbarg, mod_prop_info_t *pinfo, const char *ifname,
99 void *pval, uint_t pr_size, uint_t flags)
100 {
101     boolean_t value;
102     ill_walk_context_t ctx;
103     ill_t *ill;
104     ip_stack_t *ipst = stack->netstack_ip;
105     ip_stack_t *ipst = (ip_stack_t *)cbarg;
106     boolean_t get_def = (flags & MOD_PROP_DEFAULT);
107     boolean_t get_perm = (flags & MOD_PROP_PERM);
108     boolean_t isv6;
109     size_t nbytes = 0;
110
111     if (get_perm) {
112         nbytes = sprintf(pval, pr_size, "%d", MOD_PROP_PERM_RW);
113         goto ret;
114     } else if (get_def) {
115         nbytes = sprintf(pval, pr_size, "%d", pinfo->prop_def_bval);
116         goto ret;
117     }
118
119     /*
120      * if per interface value is not asked for return the current
121      * global value
```

```

121     */
122     if (ifname == NULL || ifname[0] == '\0') {
123         nbytes = snprintf(pval, pr_size, "%d", pinfo->prop_cur_bval);
124         goto ret;
125     }
126
127     isv6 = (pinfo->mpi_proto == MOD_PROTO_IPV6 ? B_TRUE : B_FALSE);
128     rw_enter(&ipst->ips_ill_g_lock, RW_READER);
129     if (isv6)
130         ill = ILL_START_WALK_V6(&ctx, ipst);
131     else
132         ill = ILL_START_WALK_V4(&ctx, ipst);
133     for (; ill != NULL; ill = ill_next(&ctx, ill)) {
134         /*
135          * if the property needs to be obtained on a particular
136          * interface, look for that interface.
137         */
138         if (strcmp(ifname, ill->ill_name) == 0)
139             break;
140     }
141     if (ill == NULL) {
142         rw_exit(&ipst->ips_ill_g_lock);
143         return (ENXIO);
144     }
145     value = ((ill->ill_flags & ILLF_ROUTER) ? B_TRUE : B_FALSE);
146     rw_exit(&ipst->ips_ill_g_lock);
147     nbytes = snprintf(pval, pr_size, "%d", value);
148 ret:
149     if (nbytes >= pr_size)
150         return (ENOBUFS);
151     return (0);
152 }
153 /**
154  * 'ip_debug' is a global variable. So, we will be modifying the global
155  * variable here.
156  */
157 */
158 /* ARGSUSED */
159 int
160 ip_set_debug(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
161 ip_set_debug(void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
162 const char *ifname, const void* pval, uint_t flags)
163 {
164     unsigned long new_value;
165     int err;
166
167     if (cr != NULL && secpolicy_net_config(cr, B_FALSE) != 0)
168         return (EPERM);
169
170     if ((err = mod_uint32_value(pval, pinfo, flags, &new_value)) != 0)
171         return (err);
172     ip_debug = (uint32_t)new_value;
173     return (0);
174 }
175 /**
176  * ip_debug is a global property. For default, permission and value range
177  * we retrieve the value from 'pinfo'. However for the current value we
178  * retrieve the value from the global variable 'ip_debug'
179  */
180 /* ARGSUSED */
181 int
182 ip_get_debug(netstack_t *stack, mod_prop_info_t *pinfo, const char *ifname,
183 ip_get_debug(void *cbarg, mod_prop_info_t *pinfo, const char *ifname,
184 void *pval, uint_t psiz, uint_t flags)
185 {

```

```

185     boolean_t get_def = (flags & MOD_PROP_DEFAULT);
186     boolean_t get_perm = (flags & MOD_PROP_PERM);
187     boolean_t get_range = (flags & MOD_PROP_POSSIBLE);
188     size_t nbytes;
189
190     bzero(pval, psiz);
191     if (get_perm)
192         nbytes = sprintf(pval, psiz, "%u", MOD_PROP_PERM_RW);
193     else if (get_range)
194         nbytes = sprintf(pval, psiz, "%u-%u",
195                         pinfo->prop_min_uval, pinfo->prop_max_uval);
196     else if (get_def)
197         nbytes = sprintf(pval, psiz, "%u", pinfo->prop_def_uval);
198     else
199         nbytes = sprintf(pval, psiz, "%u", ip_debug);
200     if (nbytes >= psiz)
201         return (ENOBUFS);
202     return (0);
203 }
204
205 /**
206  * Set the CGTP (multirouting) filtering status. If the status is changed
207  * from active to transparent or from transparent to active, forward the
208  * new status to the filtering module (if loaded).
209 */
210 /* ARGSUSED */
211 static int
212 ip_set_cgtp_filter(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
213 ip_set_cgtp_filter(void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
214 const char *ifname, const void* pval, uint_t flags)
215 {
216     unsigned long new_value;
217     ip_stack_t *ipst = stack->netstack_ip;
218     ip_stack_t *ipst = (ip_stack_t *)cbarg;
219     char *end;
220
221     if (flags & MOD_PROP_DEFAULT) {
222         new_value = pinfo->prop_def_bval;
223     } else {
224         if (ddi_strtoul(pval, &end, 10, &new_value) != 0 ||
225             *end != '\0' || new_value > 1) {
226             return (EINVAL);
227         }
228     }
229     if (!pinfo->prop_cur_bval && new_value) {
230         cmn_err(CE_NOTE, "IP: enabling CGTP filtering%s",
231                 ipst->ips_ip_cgtp_filter_ops == NULL ?
232                     "(module not loaded)" : "");
233     }
234     if (pinfo->prop_cur_bval && !new_value) {
235         cmn_err(CE_NOTE, "IP: disabling CGTP filtering%s",
236                 ipst->ips_ip_cgtp_filter_ops == NULL ?
237                     "(module not loaded)" : "");
238     }
239     if (ipst->ips_ip_cgtp_filter_ops != NULL) {
240         int res;
241         netstackid_t stackid = ipst->ips_netstack->netstack_stackid;
242         res = ipst->ips_ip_cgtp_filter_ops->cfo_change_state(stackid,
243                     new_value);
244         if (res)
245             return (res);
246     }
247     pinfo->prop_cur_bval = (new_value == 1 ? B_TRUE : B_FALSE);
248     ill_set_inputfn_all(ipst);
249     return (0);
250 }

```

```

249 }
250 /* Retrieve the default MTU or min-max MTU range for a given interface.
251 */
252 /* -- ill_max_frag value tells us the maximum MTU that can be handled by the
253 *   datalink. This value is advertised by the driver via DLPI messages
254 *   (DL_NOTE_SDU_SIZE/DL_INFO_ACK).
255 */
256 /* -- ill_current_frag for the most link-types will be same as ill_max_frag
257 *   to begin with. However it is dynamically computed for some link-types
258 *   like tunnels, based on the tunnel PMTU.
259 */
260 /* -- ill_mtu is the user set MTU using SIOCSLIFMTU and must lie between
261 *   (IPV6_MIN_MTU/IP_MIN_MTU) and ill_max_frag.
262 */
263 /* -- ill_user_mtu is set by in.ndpd using SIOCSLIFLINKINFO and must lie between
264 *   (IPV6_MIN_MTU/IP_MIN_MTU) and ill_max_frag.
265 */
266 */
267 */
268 int
269 ip_get_mtu(netstack_t *stack, mod_prop_info_t *pinfo, const char *ifname,
270 ip_get_mtu(void *cbarg, mod_prop_info_t *pinfo, const char *ifname,
271 void *pval, uint_t psize, uint_t flags)
272 {
273     ill_walk_context_t      ctx;
274     ill_t                   *ill;
275     ip_stack_t              *ipst = stack->netstack_ip;
276     ip_stack_t              *ipst = (ip_stack_t *)cbarg;
277     boolean_t                isv6;
278     uint32_t                 max_mtu, def_mtu;
279     size_t                  nbytes = 0;
280
281     if (!(flags & (MOD_PROP_DEFAULT|MOD_PROP_POSSIBLE)))
282         return (ENOTSUP);
283
284     if (ifname == NULL || ifname[0] == '\0')
285         return (ENOTSUP);
286
287     isv6 = (pinfo->mpi_proto == MOD_PROTO_IPV6 ? B_TRUE : B_FALSE);
288     rw_enter(&ipst->ips_ill_g_lock, RW_READER);
289     if (isv6)
290         ill = ILL_START_WALK_V6(&ctx, ipst);
291     else
292         ill = ILL_START_WALK_V4(&ctx, ipst);
293     for (; ill != NULL; ill = ill_next(&ctx, ill)) {
294         if (strcmp(ifname, ill->ill_name) == 0)
295             break;
296     }
297     if (ill == NULL) {
298         rw_exit(&ipst->ips_ill_g_lock);
299         return (ENXIO);
300     }
301     max_mtu = ill->ill_max_frag;
302     def_mtu = ill->ill_current_frag;
303     rw_exit(&ipst->ips_ill_g_lock);
304
305     if (flags & MOD_PROP_DEFAULT) {
306         nbytes = snprintf(pval, psize, "%u", def_mtu);
307     } else if (flags & MOD_PROP_POSSIBLE) {
308         uint32_t               min_mtu;
309
310         min_mtu = isv6 ? IPV6_MIN_MTU : IP_MIN_MTU;
311         nbytes = snprintf(pval, psize, "%u-%u", min_mtu, max_mtu);
312     } else {
313         return (ENOTSUP);
314     }
}

```

```

314     if (nbytes >= psize)
315         return (ENOBUFS);
316     return (0);
317 }
unchanged_portion_omitted
318 */
319 static int
320 ip_set_src_multihoming(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
321 ip_set_src_multihoming(void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
322 const char *ifname, const void* pval, uint_t flags)
323 {
324     unsigned long    new_value, old_value;
325     boolean_t        isv6;
326     ip_stack_t       *ipst = stack->netstack_ip;
327     ip_stack_t       *ipst = (ip_stack_t *)cbarg;
328     int              err;
329
330     old_value = pinfo->prop_cur_uval;
331
332     if ((err = mod_uint32_value(pval, pinfo, flags, &new_value)) != 0)
333         return (err);
334     pinfo->prop_cur_uval = new_value;
335     isv6 = (strcmp(pinfo->mpi_name, "ip6_strict_src_multihoming") == 0);
336     ip_set_src_multihoming_common(new_value, old_value, isv6, ipst);
337     return (0);
338 }
339 */
340 static int
341 ip_set_hostmodel(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
342 ip_set_hostmodel(void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
343 const char *ifname, const void* pval, uint_t flags)
344 {
345     ip_hostmodel_t   new_value, old_value;
346     ip_stack_t       *ipst = stack->netstack_ip;
347     ip_stack_t       *ipst = (ip_stack_t *)cbarg;
348     uint32_t          old_src_multihoming;
349     int              err;
350     ulong_t           tmp;
351     boolean_t         isv6;
352
353     old_value = pinfo->prop_cur_uval;
354
355     if ((err = mod_uint32_value(pval, pinfo, flags, &tmp)) != 0)
356         return (err);
357     new_value = tmp;
358     pinfo->prop_cur_uval = new_value;
359
360     switch (old_value) {
361     case IP_WEAK_ES:
362         old_src_multihoming = 0;
363         break;
364     case IP_SRC_PRI_ES:
365         old_src_multihoming = 1;
366         break;
367     case IP_STRONG_ES:
368         old_src_multihoming = 2;
369         break;
370     default:
371         ASSERT(0);
372         old_src_multihoming = IP_MAXVAL_ES;
373         break;
374     }
}

```

```

409     /*
410      * Changes to src_multihoming may require ire's to be rebound/unbound,
411      * and also require generation number resets. Changes to dst_multihoming
412      * require a simple reset of the value.
413     */
414     isv6 = (pinfo->mpi_proto == MOD_PROTO_IPV6);
415     if (new_value != old_value) {
416         switch (new_value) {
417             case IP_WEAK_ES:
418                 ip_set_src_multihoming_common(0, old_src_multihoming,
419                     isv6, ipst);
420                 if (isv6)
421                     ipst->ips_ipv6_strict_dst_multihoming = 0;
422                 else
423                     ipst->ips_ip_strict_dst_multihoming = 0;
424                 break;
425             case IP_SRC_PRI_ES:
426                 ip_set_src_multihoming_common(1, old_src_multihoming,
427                     isv6, ipst);
428                 if (isv6)
429                     ipst->ips_ipv6_strict_dst_multihoming = 0;
430                 else
431                     ipst->ips_ip_strict_dst_multihoming = 0;
432                 break;
433             case IP_STRONG_ES:
434                 ip_set_src_multihoming_common(2, old_src_multihoming,
435                     isv6, ipst);
436                 if (isv6)
437                     ipst->ips_ipv6_strict_dst_multihoming = 1;
438                 else
439                     ipst->ips_ip_strict_dst_multihoming = 1;
440                 break;
441             default:
442                 return (EINVAL);
443         }
444     }
445     return (0);
446 }

448 /* ARGSUSED */
449 int
450 ip_get_hostmodel(netstack_t *stack, mod_prop_info_t *pinfo, const char *ifname,
451 ip_get_hostmodel(void *cbarg, mod_prop_info_t *pinfo, const char *ifname,
452 void *pval, uint_t psize, uint_t flags)
453 {
454     boolean_t      isv6 = (pinfo->mpi_proto == MOD_PROTO_IPV6);
455     ip_stack_t    *ipst = stack->netstack_ip;
456     ip_stack_t    *ipst = cbarg;
457     ip_hostmodel_t hostmodel;

458     if (psize < sizeof (hostmodel))
459         return (ENOBUFS);
460     bzero(pval, psize);
461     if (!isv6) {
462         if (ipst->ips_ip_strict_src_multihoming == 0 &&
463             ipst->ips_ip_strict_dst_multihoming == 0)
464             hostmodel = IP_WEAK_ES;
465         else if (ipst->ips_ip_strict_src_multihoming == 1 &&
466             ipst->ips_ip_strict_dst_multihoming == 0)
467             hostmodel = IP_SRC_PRI_ES;
468         else if (ipst->ips_ip_strict_src_multihoming == 2 &&
469             ipst->ips_ip_strict_dst_multihoming == 1)
470             hostmodel = IP_STRONG_ES;
471         else
472             hostmodel = IP_MAXVAL_ES;
473     } else {

```

```

473         if (ipst->ips_ipv6_strict_src_multihoming == 0 &&
474             ipst->ips_ipv6_strict_dst_multihoming == 0)
475             hostmodel = IP_WEAK_ES;
476         else if (ipst->ips_ipv6_strict_src_multihoming == 1 &&
477             ipst->ips_ipv6_strict_dst_multihoming == 0)
478             hostmodel = IP_SRC_PRI_ES;
479         else if (ipst->ips_ipv6_strict_src_multihoming == 2 &&
480             ipst->ips_ipv6_strict_dst_multihoming == 1)
481             hostmodel = IP_STRONG_ES;
482         else
483             hostmodel = IP_MAXVAL_ES;
484     }
485     bcopy(&hostmodel, pval, sizeof (hostmodel));
486     return (0);
487 }



---


unchanged portion omitted

```

new/usr/src/uts/common/inet/sctp/sctp_tunables.c

10676 Mon Jul 29 18:04:20 2013

new/usr/src/uts/common/inet/sctp/sctp_tunables.c

3942 inject sanity into ipadn tcp buffer size properties

3943 _snd_lowat_fraction tcp tunable has no effect

Reviewed by: Adam Leventhal <ahl@delphix.com>

Reviewed by: Peng Dai <peng.dai@delphix.com>

```
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) 2004, 2010, Oracle and/or its affiliates. All rights reserved.
24  * Copyright (c) 2013 by Delphix. All rights reserved.
25 */
```

```
27 #include <inet/ip.h>
28 #include <inet/ip6.h>
29 #include <inet/sctp/sctp_stack.h>
30 #include <inet/sctp/sctp_impl.h>
31 #include <sys/sunddi.h>

33 /* Max size IP datagram is 64k - 1 */
34 #define SCTP_MSS_MAX_IPV4 (IP_MAXPACKET - (sizeof (iph4_t) + \
35 							   sizeof (sctp_hdr_t)))
36 #define SCTP_MSS_MAX_IPV6 (IP_MAXPACKET - (sizeof (ip6_t) + \
37 							   sizeof (sctp_hdr_t)))
38 /* Max of the above */
39 #define SCTP_MSS_MAX    SCTP_MSS_MAX_IPV4

41 /*
42  * returns the current list of listener limit configuration.
43  */
44 /* ARGSUSED */
45 static int
46 sctp_listener_conf_get(netstack_t *stack, mod_prop_info_t *pinfo,
47 	const char *ifname, void *val, uint_t psize, uint_t flags)
48 {
49 	sctp_stack_t 	*sctps = stack->netstack_sctp;
50 	sctp_stack_t 	*sctps = (sctp_stack_t *)cbarg;
51 	sctp_listener_t *sl;
52 	char 	*pval = val;
53 	size_t 	nbytes = 0, tbytes = 0;
54 	uint_t 	size;
55 	int 	err = 0;
```

1

new/usr/src/uts/common/inet/sctp/sctp_tunables.c

```
56 	bzero(pval, psize);
57 	size = psize;
58
59 	if (flags & (MOD_PROP_DEFAULT|MOD_PROP_PERM|MOD_PROP_POSSIBLE))
60 		return (0);
61
62 	mutex_enter(&sctps->sctps_listener_conf_lock);
63 	for (sl = list_head(&sctps->sctps_listener_conf); sl != NULL;
64 	sl = list_next(&sctps->sctps_listener_conf, sl)) {
65 	if (psize == size)
66 	nbytes = snprintf(pval, size, "%d:%d", sl->sl_port,
67 	sl->sl_ratio);
68 	else
69 	nbytes = snprintf(pval, size, ",%d:%d", sl->sl_port,
70 	sl->sl_ratio);
71 	size -= nbytes;
72 	pval += nbytes;
73 	tbytes += nbytes;
74 	if (tbytes >= psize) {
75 	/* Buffer overflow, stop copying information */
76 	err = ENOBUFS;
77 	break;
78  }
79 }
80 mutex_exit(&sctps->sctps_listener_conf_lock);
81 return (err);
82 }
```



```
85 /*
86  * add a new listener limit configuration.
87  */
88 /* ARGSUSED */
89 static int
90 sctp_listener_conf_add(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
91 const char *ifname, const void* pval, uint_t flags)
92 {
93 	sctp_listener_t *new_sl;
94 	sctp_listener_t *sl;
95 	long 	lport;
96 	long 	ratio;
97 	char 	*colon;
98 	sctp_stack_t 	*sctps = stack->netstack_sctp;
99 	sctp_stack_t 	*sctps = (sctp_stack_t *)cbarg;
100
101 	if (flags & MOD_PROP_DEFAULT)
102 	return (ENOTSUP);
103 	if (ddi_strtol(pval, &colon, 10, &lport) != 0 || lport <= 0 ||
104 	lport > USHRT_MAX || *colon != ':') {
105 	return (EINVAL);
106 }
107 	if (ddi_strtol(colon + 1, NULL, 10, &ratio) != 0 || ratio <= 0)
108 	return (EINVAL);
109
110 	mutex_enter(&sctps->sctps_listener_conf_lock);
111 	for (sl = list_head(&sctps->sctps_listener_conf); sl != NULL;
112 	sl = list_next(&sctps->sctps_listener_conf, sl)) {
113 	/* There is an existing entry, so update its ratio value. */
114 	if (sl->sl_port == lport) {
115 	sl->sl_ratio = ratio;
116 	mutex_exit(&sctps->sctps_listener_conf_lock);
117 	return (0);
118 }
119 }
```

2

```

121     if ((new_sl = kmalloc(sizeof (sctp_listener_t), KM_NOSLEEP)) ==
122         NULL) {
123         mutex_exit(&sctps->sctps_listener_conf_lock);
124         return (ENOMEM);
125     }
126
127     new_sl->sl_port = lport;
128     new_sl->sl_ratio = ratio;
129     list_insert_tail(&sctps->sctps_listener_conf, new_sl);
130     mutex_exit(&sctps->sctps_listener_conf_lock);
131     return (0);
132 }
133 */
134 /* remove a listener limit configuration.
135 */
136 /* ARGSUSED */
137 static int
138 sctp_listener_conf_del(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
139 sctp_listener_conf_del(void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
140     const char *ifname, const void* pval, uint_t flags)
141 {
142     sctp_listener_t *sl;
143     long lport;
144     sctp_stack_t *sctps = stack->netstack_sctp;
145     sctp_stack_t *sctps = (sctp_stack_t *)cbarg;
146
147     if (flags & MOD_PROP_DEFAULT)
148         return (ENOTSUP);
149
150     if (ddi_strtol(pval, NULL, 10, &lport) != 0 || lport <= 0 ||
151         lport > USHRT_MAX) {
152         return (EINVAL);
153     }
154     mutex_enter(&sctps->sctps_listener_conf_lock);
155     for (sl = list_head(&sctps->sctps_listener_conf); sl != NULL;
156         sl = list_next(&sctps->sctps_listener_conf, sl)) {
157         if (sl->sl_port == lport) {
158             list_remove(&sctps->sctps_listener_conf, sl);
159             mutex_exit(&sctps->sctps_listener_conf_lock);
160             kmem_free(sl, sizeof (sctp_listener_t));
161             return (0);
162         }
163     }
164     mutex_exit(&sctps->sctps_listener_conf_lock);
165 }
166 static int
167 sctp_set_buf_prop(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
168     const char *ifname, const void *pval, uint_t flags)
169 {
170     return (mod_set_buf_prop(stack->netstack_sctp->sctps_propinfo_tbl,
171         stack, cr, pinfo, ifname, pval, flags));
172 }
173
174 static int
175 sctp_get_buf_prop(netstack_t *stack, mod_prop_info_t *pinfo, const char *ifname,
176     void *val, uint_t psize, uint_t flags)
177 {
178     return (mod_get_buf_prop(stack->netstack_sctp->sctps_propinfo_tbl,
179         stack, pinfo, ifname, val, psize, flags));
180 }
181 */
182

```

```

184     * All of these are alterable, within the min/max values given, at run time.
185     *
186     * Note: All those tunables which do not start with "_" are Committed and
187     * therefore are public. See PSARC 2010/080.
188     */
189 mod_prop_info_t sctp_propinfo_tbl[] = {
190     { "_max_init_retr", MOD_PROTO_SCTP,
191         mod_set_uint32, mod_get_uint32,
192         {0, 128, 8}, {8} },
193     { "pa_max_retr", MOD_PROTO_SCTP,
194         mod_set_uint32, mod_get_uint32,
195         {1, 128, 10}, {10} },
196     { "pp_max_retr", MOD_PROTO_SCTP,
197         mod_set_uint32, mod_get_uint32,
198         {1, 128, 5}, {5} },
199     { "cwnd_max", MOD_PROTO_SCTP,
200         mod_set_uint32, mod_get_uint32,
201         {128, ULP_MAX_BUF, 1024*1024}, {1024*1024} },
202         {128, (1<<30), 1024*1024}, {1024*1024} },
203     { "smallest_nonpriv_port", MOD_PROTO_SCTP,
204         mod_set_uint32, mod_get_uint32,
205         {1024, (32*1024), 1024}, {1024} },
206     { "ipv4_ttl", MOD_PROTO_SCTP,
207         mod_set_uint32, mod_get_uint32,
208         {1, 255, 64}, {64} },
209     { "heartbeat_interval", MOD_PROTO_SCTP,
210         mod_set_uint32, mod_get_uint32,
211         {0, 1*DAYS, 30*SECONDS}, {30*SECONDS} },
212     { "initial_mtu", MOD_PROTO_SCTP,
213         mod_set_uint32, mod_get_uint32,
214         {68, 65535, 1500}, {1500} },
215     { "mtu_probe_interval", MOD_PROTO_SCTP,
216         mod_set_uint32, mod_get_uint32,
217         {0, 1*DAYS, 10*MINUTES}, {10*MINUTES} },
218     { "new_secret_interval", MOD_PROTO_SCTP,
219         mod_set_uint32, mod_get_uint32,
220         {0, 1*DAYS, 2*MINUTES}, {2*MINUTES} },
221     /* tunable - 10 */
222     { "deferred_ack_interval", MOD_PROTO_SCTP,
223         mod_set_uint32, mod_get_uint32,
224         {10*MS, 1*MINUTES, 100*MS}, {100*MS} },
225     { "snd_lowat_fraction", MOD_PROTO_SCTP,
226         mod_set_uint32, mod_get_uint32,
227         {0, 16, 0}, {0} },
228     { "ignore_path_mtu", MOD_PROTO_SCTP,
229         mod_set_boolean, mod_get_boolean,
230         {B_FALSE}, {B_FALSE} },
231     { "initial_ssthresh", MOD_PROTO_SCTP,
232         mod_set_uint32, mod_get_uint32,
233         {1024, UINT32_MAX, SCTP_RECV_HIWATER}, {SCTP_RECV_HIWATER} },
234     { "smallest_anon_port", MOD_PROTO_SCTP,
235         mod_set_uint32, mod_get_uint32,
236         {0, 1024, 1024}, {1024} },
237     { "initial_cwnd", MOD_PROTO_SCTP,
238         mod_set_uint32, mod_get_uint32,
239         {0, 1024, 1024}, {1024} },
240     { "initial_retrans", MOD_PROTO_SCTP,
241         mod_set_uint32, mod_get_uint32,
242         {0, 1024, 1024}, {1024} },
243     { "initial_mtu", MOD_PROTO_SCTP,
244         mod_set_uint32, mod_get_uint32,
245         {0, 1024, 1024}, {1024} },
246     { "initial_cwnd", MOD_PROTO_SCTP,
247         mod_set_uint32, mod_get_uint32,
248         {0, 1024, 1024}, {1024} },
249     { "initial_retrans", MOD_PROTO_SCTP,
250         mod_set_uint32, mod_get_uint32,
251         {0, 1024, 1024}, {1024} },
252     { "initial_mtu", MOD_PROTO_SCTP,
253         mod_set_uint32, mod_get_uint32,
254         {0, 1024, 1024}, {1024} },
255     { "initial_cwnd", MOD_PROTO_SCTP,
256         mod_set_uint32, mod_get_uint32,
257         {0, 1024, 1024}, {1024} },
258     { "initial_retrans", MOD_PROTO_SCTP,
259         mod_set_uint32, mod_get_uint32,
260         {0, 1024, 1024}, {1024} },
261     { "initial_mtu", MOD_PROTO_SCTP,
262         mod_set_uint32, mod_get_uint32,
263         {0, 1024, 1024}, {1024} },
264     { "initial_cwnd", MOD_PROTO_SCTP,
265         mod_set_uint32, mod_get_uint32,
266         {0, 1024, 1024}, {1024} },
267     { "initial_retrans", MOD_PROTO_SCTP,
268         mod_set_uint32, mod_get_uint32,
269         {0, 1024, 1024}, {1024} },
270     { "initial_mtu", MOD_PROTO_SCTP,
271         mod_set_uint32, mod_get_uint32,
272         {0, 1024, 1024}, {1024} },
273     { "initial_cwnd", MOD_PROTO_SCTP,
274         mod_set_uint32, mod_get_uint32,
275         {0, 1024, 1024}, {1024} },
276     { "initial_retrans", MOD_PROTO_SCTP,
277         mod_set_uint32, mod_get_uint32,
278         {0, 1024, 1024}, {1024} },
279     { "initial_mtu", MOD_PROTO_SCTP,
280         mod_set_uint32, mod_get_uint32,
281         {0, 1024, 1024}, {1024} },
282     { "initial_cwnd", MOD_PROTO_SCTP,
283         mod_set_uint32, mod_get_uint32,
284         {0, 1024, 1024}, {1024} },
285     { "initial_retrans", MOD_PROTO_SCTP,
286         mod_set_uint32, mod_get_uint32,
287         {0, 1024, 1024}, {1024} },
288 }
289

```

```

249           {1024, ULP_MAX_PORT, 32*1024}, {32*1024} },
250
251     { "largest_anon_port", MOD_PROTO_SCTP,
252       mod_set_uint32, mod_get_uint32,
253       {1024, ULP_MAX_PORT, ULP_MAX_PORT}, {ULP_MAX_PORT} },
254
255     { "send_buf", MOD_PROTO_SCTP,
256       sctp_set_buf_prop, sctp_get_buf_prop,
257       {SCTP_XMIT_LOWATER, ULP_MAX_BUF, SCTP_XMIT_HIWATER},
258
259     { "send_maxbuf", MOD_PROTO_SCTP,
260       mod_set_uint32, mod_get_uint32,
261       {SCTP_XMIT_LOWATER, (1<<30), SCTP_XMIT_HIWATER},
262       {SCTP_XMIT_HIWATER} },
263
264     { "xmit_lowat", MOD_PROTO_SCTP,
265       mod_set_uint32, mod_get_uint32,
266       {SCTP_XMIT_LOWATER, ULP_MAX_BUF, SCTP_XMIT_LOWATER},
267       {SCTP_XMIT_LOWATER, (1<<30), SCTP_XMIT_LOWATER},
268       {SCTP_XMIT_LOWATER} },
269
270     { "recv_buf", MOD_PROTO_SCTP,
271       sctp_set_buf_prop, sctp_get_buf_prop,
272       {SCTP_RECV_LOWATER, ULP_MAX_BUF, SCTP_RECV_HIWATER},
273
274     /* tunable - 20 */
275     { "_rtt_updates", MOD_PROTO_SCTP,
276       mod_set_uint32, mod_get_uint32,
277       {0, 65536, 20}, {20} },
278
279     { "_ipv6_hoplimit", MOD_PROTO_SCTP,
280       mod_set_uint32, mod_get_uint32,
281       {0, IPV6_MAX_HOPS, IPV6_DEFAULT_HOPS}, {IPV6_DEFAULT_HOPS} },
282
283     { "_rto_min", MOD_PROTO_SCTP,
284       mod_set_uint32, mod_get_uint32,
285       {500*MS, 60*SECONDS, 1*SECONDS}, {1*SECONDS} },
286
287     { "_rto_max", MOD_PROTO_SCTP,
288       mod_set_uint32, mod_get_uint32,
289       {1*SECONDS, 60000*SECONDS, 60*SECONDS}, {60*SECONDS} },
290
291     { "_rto_initial", MOD_PROTO_SCTP,
292       mod_set_uint32, mod_get_uint32,
293       {1*SECONDS, 60000*SECONDS, 3*SECONDS}, {3*SECONDS} },
294
295     { "_cookie_life", MOD_PROTO_SCTP,
296       mod_set_uint32, mod_get_uint32,
297       {10*MS, 60000*SECONDS, 60*SECONDS}, {60*SECONDS} },
298
299     { "_max_in_streams", MOD_PROTO_SCTP,
300       mod_set_uint32, mod_get_uint32,
301       {1, UINT16_MAX, 32}, {32} },
302
303     { "_initial_out_streams", MOD_PROTO_SCTP,
304       mod_set_uint32, mod_get_uint32,
305       {1, UINT16_MAX, 32}, {32} },

```

```

307           {"_shutack_wait_bound", MOD_PROTO_SCTP,
308             mod_set_uint32, mod_get_uint32,
309             {0, 300*SECONDS, 60*SECONDS}, {60*SECONDS} },
310
311           {"_maxburst", MOD_PROTO_SCTP,
312             mod_set_uint32, mod_get_uint32,
313             {2, 8, 4}, {4} },
314
315           /* tunable - 30 */
316           {"_addip_enabled", MOD_PROTO_SCTP,
317             mod_set_boolean, mod_get_boolean,
318             {B_FALSE}, {B_FALSE} },
319
320           {"_recv_hiwat_minmss", MOD_PROTO_SCTP,
321             mod_set_uint32, mod_get_uint32,
322             {1, 65536, 4}, {4} },
323
324           {"_slow_start_initial", MOD_PROTO_SCTP,
325             mod_set_uint32, mod_get_uint32,
326             {1, 16, 4}, {4} },
327
328           {"_slow_start_after_idle", MOD_PROTO_SCTP,
329             mod_set_uint32, mod_get_uint32,
330             {1, 16384, 4}, {4} },
331
332           {"_prscpt_enabled", MOD_PROTO_SCTP,
333             mod_set_boolean, mod_get_boolean,
334             {B_TRUE}, {B_TRUE} },
335
336           {"_fast_rxt_thresh", MOD_PROTO_SCTP,
337             mod_set_uint32, mod_get_uint32,
338             {1, 10000, 3}, {3} },
339
340           {"_deferred_acks_max", MOD_PROTO_SCTP,
341             mod_set_uint32, mod_get_uint32,
342             {1, 16, 2}, {2} },
343
344           /*
345            * sctp_wroff_xtra is the extra space in front of SCTP/IP header
346            * for link layer header. It has to be a multiple of 8.
347            */
348           {"_wroff_xtra", MOD_PROTO_SCTP,
349             mod_set_aligned, mod_get_uint32,
350             {0, 256, 32}, {32} },
351
352           {"_extra_priv_ports", MOD_PROTO_SCTP,
353             mod_set_extra_privports, mod_get_extra_privports,
354             {1, ULP_MAX_PORT, 0}, {0} },
355
356           {"_listener_limit_conf", MOD_PROTO_SCTP,
357             NULL, sctp_listener_conf_get, {0}, {0} },
358
359           {"_listener_limit_conf_add", MOD_PROTO_SCTP,
360             sctp_listener_conf_add, NULL, {0}, {0} },
361
362           {"_listener_limit_conf_del", MOD_PROTO_SCTP,
363             sctp_listener_conf_del, NULL, {0}, {0} },
364
365           {"?", MOD_PROTO_SCTP, NULL, mod_get_allprop, {0}, {0} },
366           {NULL, 0, NULL, NULL, {0}, {0} }
367
368 };

```

unchanged portion omitted

new/usr/src/uts/common/inet/tcp/tcp.c

129719 Mon Jul 29 18:04:21 2013

new/usr/src/uts/common/inet/tcp/tcp.c

3942 inject sanity into ipadmin tcp buffer size properties

3943 _snd_lowat_fraction tcp tunable has no effect

Reviewed by: Adam Leventhal <ahl@delphix.com>

Reviewed by: Peng Dai <peng.dai@delphix.com>

```
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) 2011, Joyent Inc. All rights reserved.
25 * Copyright (c) 2011 Nexenta Systems, Inc. All rights reserved.
26 * Copyright (c) 2013 by Delphix. All rights reserved.
27 */
28 /* Copyright (c) 1990 Mentor Inc. */
```

```
30 #include <sys/types.h>
31 #include <sys/stream.h>
32 #include <sys/strsun.h>
33 #include <sys/strsubr.h>
34 #include <sys/stropts.h>
35 #include <sys/strlog.h>
36 #define _SUN_TPI_VERSION 2
37 #include <sys/tihdr.h>
38 #include <sys/timod.h>
39 #include <sys/ddi.h>
40 #include <sys/sunddi.h>
41 #include <sys/suntpi.h>
42 #include <sys/xti_inet.h>
43 #include <sys/cmn_err.h>
44 #include <sys/debug.h>
45 #include <sys/sdt.h>
46 #include <sys/vtrace.h>
47 #include <sys/kmem.h>
48 #include <sys/ethernet.h>
49 #include <sys/cpuvar.h>
50 #include <sys/dlpi.h>
51 #include <sys/pattr.h>
52 #include <sys/policy.h>
53 #include <sys/priv.h>
54 #include <sys/zone.h>
55 #include <sys/sunlidi.h>
```

```
57 #include <sys/errno.h>
58 #include <sys/signal.h>
```

1

new/usr/src/uts/common/inet/tcp/tcp.c

```
59 #include <sys/socket.h>
60 #include <sys/socketvar.h>
61 #include <sys/sockio.h>
62 #include <sys/isa_defs.h>
63 #include <sys/md5.h>
64 #include <sys/random.h>
65 #include <sys/uio.h>
66 #include <sys/sysm.h>
67 #include <netinet/in.h>
68 #include <netinet/tcp.h>
69 #include <netinet/ip6.h>
70 #include <netinet/icmp6.h>
71 #include <net/if.h>
72 #include <net/route.h>
73 #include <inet/ipsec_impl.h>

75 #include <inet/common.h>
76 #include <inet/ip.h>
77 #include <inet/ip_impl.h>
78 #include <inet/ip6.h>
79 #include <inet/ip_ndp.h>
80 #include <inet/proto_set.h>
81 #include <inet/mib2.h>
82 #include <inet/optcom.h>
83 #include <inet/snmpcom.h>
84 #include <inet/kstatcom.h>
85 #include <inet/tcp.h>
86 #include <inet/tcp_impl.h>
87 #include <inet/tcp_cluster.h>
88 #include <inet/udp_impl.h>
89 #include <net/pfkeyv2.h>
90 #include <inet/ipdrop.h>

92 #include <inet/ipclassifier.h>
93 #include <inet/ip_ire.h>
94 #include <inet/ip_ftable.h>
95 #include <inet/ip_if.h>
96 #include <inet/ipp_common.h>
97 #include <inet/ip_rts.h>
98 #include <inet/ip_netinfo.h>
99 #include <sys/squeue_impl.h>
100 #include <sys/squeue.h>
101 #include <sys/tsol/label.h>
102 #include <sys/tsol/tnet.h>
103 #include <rpc/pmap_prot.h>
104 #include <sys/callo.h>

106 /*
107 * TCP Notes: aka FireEngine Phase I (PSARC 2002/433)
108 *
109 * (Read the detailed design doc in PSARC case directory)
110 *
111 * The entire tcp state is contained in tcp_t and conn_t structure
112 * which are allocated in tandem using ipcl_conn_create() and passing
113 * IPCL_TCPCONN as a flag. We use 'conn_ref' and 'conn_lock' to protect
114 * the references on the tcp_t. The tcp_t structure is never compressed
115 * and packets always land on the correct TCP perimeter from the time
116 * eager is created till the time tcp_t dies (as such the old mentat
117 * TCP global queue is not used for detached state and no IPSEC checking
118 * is required). The global queue is still allocated to send out resets
119 * for connection which have no listeners and IP directly calls
120 * tcp_xmit_listeners_reset() which does any policy check.
121 *
122 * Protection and Synchronisation mechanism:
123 *
124 * The tcp data structure does not use any kind of lock for protecting
```

2

```

125 * its state but instead uses 'squeues' for mutual exclusion from various
126 * read and write side threads. To access a tcp member, the thread should
127 * always be behind squeue (via squeue_enter with flags as SQ_FILL, SQ_PROCESS,
128 * or SQ_NODRAIN). Since the squeues allow a direct function call, caller
129 * can pass any tcp function having prototype of edesc_t as argument
130 * (different from traditional STREAMs model where packets come in only
131 * designated entry points). The list of functions that can be directly
132 * called via squeue are listed before the usual function prototype.
133 *
134 * Referencing:
135 *
136 * TCP is MT-Hot and we use a reference based scheme to make sure that the
137 * tcp structure doesn't disappear when its needed. When the application
138 * creates an outgoing connection or accepts an incoming connection, we
139 * start out with 2 references on 'conn_ref'. One for TCP and one for IP.
140 * The IP reference is just a symbolic reference since ip_tcpclose()
141 * looks at tcp structure after tcp_close_output() returns which could
142 * have dropped the last TCP reference. So as long as the connection is
143 * in attached state i.e. !TCP_IS_DETACHED, we have 2 references on the
144 * conn_t. The classifier puts its own reference when the connection is
145 * inserted in listen or connected hash. Anytime a thread needs to enter
146 * the tcp connection perimeter, it retrieves the conn/tcp from q->ptr
147 * on write side or by doing a classify on read side and then puts a
148 * reference on the conn before doing squeue_enter/tryenter/fill. For
149 * read side, the classifier itself puts the reference under fanout lock
150 * to make sure that tcp can't disappear before it gets processed. The
151 * squeue will drop this reference automatically so the called function
152 * doesn't have to do a DEC_REF.
153 *
154 * Opening a new connection:
155 *
156 * The outgoing connection open is pretty simple. tcp_open() does the
157 * work in creating the conn/tcp structure and initializing it. The
158 * squeue assignment is done based on the CPU the application
159 * is running on. So for outbound connections, processing is always done
160 * on application CPU which might be different from the incoming CPU
161 * being interrupted by the NIC. An optimal way would be to figure out
162 * the NIC <-> CPU binding at listen time, and assign the outgoing
163 * connection to the squeue attached to the CPU that will be interrupted
164 * for incoming packets (we know the NIC based on the bind IP address).
165 * This might seem like a problem if more data is going out but the
166 * fact is that in most cases the transmit is ACK driven transmit where
167 * the outgoing data normally sits on TCP's xmit queue waiting to be
168 * transmitted.
169 *
170 * Accepting a connection:
171 *
172 * This is a more interesting case because of various races involved in
173 * establishing a eager in its own perimeter. Read the meta comment on
174 * top of tcp_input_listener(). But briefly, the squeue is picked by
175 * ip_fanout based on the ring or the sender (if loopback).
176 *
177 * Closing a connection:
178 *
179 * The close is fairly straight forward. tcp_close() calls tcp_close_output()
180 * via squeue to do the close and mark the tcp as detached if the connection
181 * was in state TCPS_ESTABLISHED or greater. In the later case, TCP keep its
182 * reference but tcp_close() drop IP's reference always. So if tcp was
183 * not killed, it is sitting in time_wait list with 2 reference - 1 for TCP
184 * and 1 because it is in classifier's connected hash. This is the condition
185 * we use to determine that its OK to clean up the tcp outside of squeue
186 * when time wait expires (check the ref under fanout and conn_lock and
187 * if it is 2, remove it from fanout hash and kill it).
188 *
189 * Although close just drops the necessary references and marks the
190 * tcp_detached state, tcp_close needs to know the tcp_detached has been

```

```

191 * set (under squeue) before letting the STREAM go away (because a
192 * inbound packet might attempt to go up the STREAM while the close
193 * has happened and tcp_detached is not set). So a special lock and
194 * flag is used along with a condition variable (tcp_closelock, tcp_closed,
195 * and tcp_closecv) to signal tcp_close that tcp_close_out() has marked
196 * tcp_detached.
197 *
198 * Special provisions and fast paths:
199 *
200 * We make special provisions for sockfs by marking tcp_issocket
201 * whenever we have only sockfs on top of TCP. This allows us to skip
202 * putting the tcp in acceptor hash since a sockfs listener can never
203 * become acceptor and also avoid allocating a tcp_t for acceptor STREAM
204 * since eager has already been allocated and the accept now happens
205 * on acceptor STREAM. There is a big blob of comment on top of
206 * tcp_input_listener explaining the new accept. When socket is POP'd,
207 * sockfs sends us an ioctl to mark the fact and we go back to old
208 * behaviour. Once tcp_issocket is unset, its never set for the
209 * life of that connection.
210 *
211 * IPsec notes :
212 *
213 * Since a packet is always executed on the correct TCP perimeter
214 * all IPsec processing is defered to IP including checking new
215 * connections and setting IPSEC policies for new connection. The
216 * only exception is tcp_xmit_listeners_reset() which is called
217 * directly from IP and needs to policy check to see if TH_RST
218 * can be sent out.
219 */
220 */
221 * Values for squeue switch:
222 * 1: SQ_NODRAIN
223 * 2: SQ_PROCESS
224 * 3: SQ_FILL
225 */
226 int tcp_squeue_wput = 2; /* /etc/systems */
227 int tcp_squeue_flag;
228
229 /* To prevent memory hog, limit the number of entries in tcp_free_list
230 * to 1% of available memory / number of cpus
231 */
232 uint_t tcp_free_list_max_cnt = 0;
233
234 #define TCP_XMIT_LOWATER 4096
235 #define TCP_XMIT_HIWATER 49152
236 #define TCP_RECV_LOWATER 2048
237 #define TCP_RECV_HIWATER 128000
238
239 /* transport interface data unit size */
240 /*
241 #define TCP_ACCEPTOR_FANOUT_SIZE 512
242
243 #ifdef _ILP32
244 #define TCP_ACCEPTOR_HASH(accid)
245 ((uint_t)(accid) >> 8) & (TCP_ACCEPTOR_FANOUT_SIZE - 1))
246 #else
247 #define TCP_ACCEPTOR_HASH(accid)
248 ((uint_t)(accid) & (TCP_ACCEPTOR_FANOUT_SIZE - 1))
249 #endif /* _ILP32 */
250 */

```

```

252 * Minimum number of connections which can be created per listener. Used
253 * when the listener connection count is in effect.
254 */
255 static uint32_t tcp_min_conn_listener = 2;

257 uint32_t tcp_early_abort = 30;

259 /* TCP Timer control structure */
260 typedef struct tcpt_s {
261     pfv_t      tcpt_pfv;          /* The routine we are to call */
262     tcpt_t     *tcpt_tcp;         /* The parameter we are to pass in */
263 } tcpt_t;
unchanged portion omitted

2590 conn_t *
2591 tcp_create_common(cred_t *credp, boolean_t isv6, boolean_t issocket,
2592                     int *errorp)
2593 {
2594     tcp_t          *tcp = NULL;
2595     conn_t         *connp;
2596     zoneid_t       zoneid;
2597     tcp_stack_t    *tcps;
2598     squeue_t       *sqp;

2600     ASSERT(errorp != NULL);
2601     /*
2602     * Find the proper zoneid and netstack.
2603     */
2604     /*
2605     * Special case for install: miniroot needs to be able to
2606     * access files via NFS as though it were always in the
2607     * global zone.
2608     */
2609     if (credp == kcred && nfs_global_client_only != 0) {
2610         zoneid = GLOBAL_ZONEID;
2611         tcps = netstack_find_by_stackid(GLOBAL_NETSTACKID)->
2612               netstack_tcp;
2613         ASSERT(tcps != NULL);
2614     } else {
2615         netstack_t *ns;
2616         int err;

2617         if ((err = secpolicy_basic_net_access(credp)) != 0) {
2618             *errorp = err;
2619             return (NULL);
2620         }

2621         ns = netstack_find_by_cred(credp);
2622         ASSERT(ns != NULL);
2623         tcps = ns->netstack_tcp;
2624         ASSERT(tcps != NULL);

2625         /*
2626         * For exclusive stacks we set the zoneid to zero
2627         * to make TCP operate as if in the global zone.
2628         */
2629         if (tcps->tcps_netstack->netstack_stackid !=
2630             GLOBAL_NETSTACKID)
2631             zoneid = GLOBAL_ZONEID;
2632         else
2633             zoneid = crgetzoneid(credp);
2634     }

2635     sqp = IP_SQUEUE_GET((uint_t)gethrtime());
2636     connp = (conn_t *)tcp_get_conn(sqp, tcps);
2637     /*

```

```

2642             * Both tcp_get_conn and netstack_find_by_cred incremented refcnt,
2643             * so we drop it by one.
2644             */
2645             netstack_rele(tcps->tcps_netstack);
2646             if (connp == NULL) {
2647                 *errorp = ENOSR;
2648                 return (NULL);
2649             }
2650             ASSERT(connp->conn_ixa->ixa_protocol == connp->conn_proto);
2651             connp->conn_sqp = sqp;
2652             connp->conn_initial_sqp = connp->conn_sqp;
2653             connp->conn_ixa->ixa_sqp = connp->conn_sqp;
2654             tcp = connp->conn_tcp;
2655             */

2656             /*
2657             * Besides asking IP to set the checksum for us, have conn_ip_output
2658             * to do the following checks when necessary:
2659             *
2660             * IXAF_VERIFY_SOURCE: drop packets when our outer source goes invalid
2661             * IXAF_VERIFY_PMTU: verify PMTU changes
2662             * IXAF_VERIFY_LSO: verify LSO capability changes
2663             */
2664             connp->conn_ixa->ixa_flags |= IXAF_SET_ULP_CKSUM | IXAF_VERIFY_SOURCE |
2665                 IXAF_VERIFY_PMTU | IXAF_VERIFY_LSO;
2666             if (!tcps->tcps_dev_flow_ctl)
2667                 connp->conn_ixa->ixa_flags |= IXAF_NO_DEV_FLOW_CTL;

2668             if (isv6) {
2669                 connp->conn_ixa->ixa_src_preferences = IPV6_PREFER_SRC_DEFAULT;
2670                 connp->conn_ipversion = IPV6_VERSION;
2671                 connp->conn_family = AF_INET6;
2672                 tcp->tcp_mss = tcps->tcps_mss_def_ipv6;
2673                 connp->conn_default_ttl = tcps->tcps_ipv6_hoplimit;
2674             } else {
2675                 connp->conn_ipversion = IPV4_VERSION;
2676                 connp->conn_family = AF_INET;
2677                 tcp->tcp_mss = tcps->tcps_mss_def_ipv4;
2678                 connp->conn_default_ttl = tcps->tcps_ipv4_ttl;
2679             }
2680             connp->conn_xmit_ipp.ipp_unicast_hops = connp->conn_default_ttl;
2681             crhold(credp);
2682             connp->conn_cred = credp;
2683             connp->conn_cpid = curproc->p_pid;
2684             connp->conn_open_time = ddi_get_lbolt64();

2685             /*
2686             * Cache things in the ixa without any refhold */
2687             ASSERT(!(connp->conn_ixa->ixa_free_flags & IXA_FREE_CRED));
2688             connp->conn_ixa->ixa_cred = credp;
2689             connp->conn_ixa->ixa_cpid = connp->conn_cpid;
2690             connp->conn_zoneid = zoneid;
2691             /* conn_allzones can not be set this early, hence no IPCL_ZONEID */
2692             connp->conn_ixa->ixa_zoneid = zoneid;
2693             connp->conn_mlp_type = mlptSingle;
2694             ASSERT(connp->conn_netstack == tcps->tcps_netstack);
2695             ASSERT(tcp->tcp_tcps == tcps);

2696             /*
2697             * If the caller has the process-wide flag set, then default to MAC
2698             * exempt mode. This allows read-down to unlabeled hosts.
2699             */
2700             if (getpflags(NET_MAC_AWARE, credp) != 0)
2701                 connp->conn_mac_mode = CONN_MAC_AWARE;

```

```
2709     connp->conn_zone_is_global = (crgetzoneid(credp) == GLOBAL_ZONEID);
2711     if (issocket) {
2712         tcp->tcp_issocket = 1;
2713     }
2715     connp->conn_rcvbuf = tcps->tcps_recv_hiwat;
2716     connp->conn_sndbuf = tcps->tcps_xmit_hiwat;
2717     if (tcps->tcps_snd_lowat_fraction != 0) {
2718         connp->conn_sndlowat = connp->conn_sndbuf /
2719             tcps->tcps_snd_lowat_fraction;
2720     } else {
2721         connp->conn_sndlowat = tcps->tcps_xmit_lowat;
2722     }
2723     connp->conn_so_type = SOCK_STREAM;
2724     connp->conn_wroff = connp->conn_ht_iphc_allocated +
2725         tcps->tcps_wroff_xtra;
2727     SOCK_CONNID_INIT(tcp->tcp_connid);
2728     /* DTrace ignores this - it isn't a tcp:::state-change */
2729     tcp->tcp_state = TCPS_IDLE;
2730     tcp_init_values(tcp, NULL);
2731     return (connp);
2732 }
```

unchanged portion omitted

```
new/usr/src/uts/common/inet/tcp/tcp_tunables.c
```

```
1
```

```
*****
15246 Mon Jul 29 18:04:23 2013
new/usr/src/uts/common/inet/tcp/tcp_tunables.c
3942 inject sanity into ipadm tcp buffer size properties
3943 _snd_lowat_fraction tcp tunable has no effect
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Peng Dai <peng.dai@delphix.com>
*****
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) 1991, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 2011, Joyent Inc. All rights reserved.
24 * Copyright 2013 Nexenta Systems, Inc. All rights reserved.
25 * Copyright (c) 2013 by Delphix. All rights reserved.
26 */
27 /* Copyright (c) 1990 Mentor Inc. */

29 #include <inet/ip.h>
30 #include <inet/tcp_impl.h>
31 #include <sys/multidata.h>
32 #include <sys/sunddi.h>

34 /* Max size IP datagram is 64k - 1 */
35 #define TCP_MSS_MAX_IPV4 (IP_MAXPACKET - (sizeof (iph4_t) + sizeof (tcp4h_t)))
36 #define TCP_MSS_MAX_IPV6 (IP_MAXPACKET - (sizeof (ip6_t) + sizeof (tcp6h_t)))

38 /* Max of the above */
39 #define TCP_MSS_MAX TCP_MSS_MAX_IPV4

40 #define TCP_XMIT_LOWATER 4096
41 #define TCP_XMIT_HIWATER 49152
42 #define TCP_RECV_LOWATER 2048
43 #define TCP_RECV_HIWATER 128000

41 /*
42 * Set the RFC 1948 pass phrase
43 */
44 /* ARGSUSED */
45 static int
46 tcp_set_1948phrase(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
47 tcp_set_1948phrase(void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
48 const char *ifname, const void* pr_val, uint_t flags)
49 {
50     tcp_stack_t *tcps = (tcp_stack_t *)cbarg;
51     if (flags & MOD_PROP_DEFAULT)
52         return (ENOTSUP);
```

```
new/usr/src/uts/common/inet/tcp/tcp_tunables.c
```

```
2
```

```
52     /*
53      * Basically, value contains a new pass phrase. Pass it along!
54      */
55     tcp_iss_key_init((uint8_t *)pr_val, strlen(pr_val),
56                      stack->netstack_tcp);
57     tcp_iss_key_init((uint8_t *)pr_val, strlen(pr_val), tcps);
58 }

59 /*
60  * returns the current list of listener limit configuration.
61 */
62 /* ARGSUSED */
63 static int
64 tcp_listener_conf_get(netstack_t *stack, mod_prop_info_t *pinfo,
65 const char *ifname, void *val, uint_t psize, uint_t flags)
66 tcp_listener_conf_get(void *cbarg, mod_prop_info_t *pinfo, const char *ifname,
67 void *val, uint_t psize, uint_t flags)
68 {
69     tcp_stack_t *tcps = stack->netstack_tcp;
70     tcp_stack_t *tcps = (tcp_stack_t *)cbarg;
71     tcp_listener_t *tl;
72     char *pval = val;
73     size_t nbytes = 0, tbytes = 0;
74     uint_t size;
75     int err = 0;

76     bzero(pval, psize);
77     size = psize;

78     if (flags & (MOD_PROP_DEFAULT|MOD_PROP_PERM|MOD_PROP_POSSIBLE))
79         return (0);

80     mutex_enter(&tcps->tcps_listener_conf_lock);
81     for (tl = list_head(&tcps->tcps_listener_conf); tl != NULL;
82          tl = list_next(&tcps->tcps_listener_conf, tl)) {
83         if (psize == size)
84             nbytes = snprintf(pval, size, "%d:%d", tl->tl_port,
85                               tl->tl_ratio);
86         else
87             nbytes = snprintf(pval, size, ",%d:%d", tl->tl_port,
88                               tl->tl_ratio);
88         size -= nbytes;
89         pval += nbytes;
90         tbytes += nbytes;
91         if (tbytes >= psize) {
92             /* Buffer overflow, stop copying information */
93             err = ENOBUFS;
94             break;
95         }
96     }
97 }

98 }

100    mutex_exit(&tcps->tcps_listener_conf_lock);
101    return (err);
102 }

104 /*
105  * add a new listener limit configuration.
106 */
107 /* ARGSUSED */
108 static int
109 tcp_listener_conf_add(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
110 tcp_listener_conf_add(void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
111 const char *ifname, const void* pval, uint_t flags)
```

```

112     tcp_listener_t *new_tl;
113     tcp_listener_t *tl;
114     long lport;
115     long ratio;
116     char *colon;
117     tcp_stack_t *tcps = stack->netstack_tcp;
118     tcp_stack_t *tcps = (tcp_stack_t *)cbarg;
119
120     if (flags & MOD_PROP_DEFAULT)
121         return (ENOTSUP);
122
123     if (ddi_strtol(pval, &colon, 10, &lport) != 0 || lport <= 0 ||
124         lport > USHRT_MAX || *colon != ':') {
125         return (EINVAL);
126     }
127     if (ddi_strtol(colon + 1, NULL, 10, &ratio) != 0 || ratio <= 0)
128         return (EINVAL);
129
130     mutex_enter(&tcps->tcps_listener_conf_lock);
131     for (tl = list_head(&tcps->tcps_listener_conf); tl != NULL;
132          tl = list_next(&tcps->tcps_listener_conf, tl)) {
133         /* There is an existing entry, so update its ratio value. */
134         if (tl->tl_port == lport) {
135             tl->tl_ratio = ratio;
136             mutex_exit(&tcps->tcps_listener_conf_lock);
137             return (0);
138         }
139     }
140
141     if ((new_tl = kmem_alloc(sizeof (tcp_listener_t), KM_NOSLEEP)) ==
142         NULL) {
143         mutex_exit(&tcps->tcps_listener_conf_lock);
144         return (ENOMEM);
145     }
146
147     new_tl->tl_port = lport;
148     new_tl->tl_ratio = ratio;
149     list_insert_tail(&tcps->tcps_listener_conf, new_tl);
150     mutex_exit(&tcps->tcps_listener_conf_lock);
151 }
152
153 /* remove a listener limit configuration.
154 */
155 /* ARGSUSED */
156 static int
157 tcp_listener_conf_del(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
158 void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
159 const char *ifname, const void* pval, uint_t flags)
160 {
161     tcp_listener_t *tl;
162     long lport;
163     tcp_stack_t *tcps = stack->netstack_tcp;
164     tcp_stack_t *tcps = (tcp_stack_t *)cbarg;
165
166     if (flags & MOD_PROP_DEFAULT)
167         return (ENOTSUP);
168
169     if (ddi_strtol(pval, NULL, 10, &lport) != 0 || lport <= 0 ||
170         lport > USHRT_MAX) {
171         return (EINVAL);
172     }
173     mutex_enter(&tcps->tcps_listener_conf_lock);
174     for (tl = list_head(&tcps->tcps_listener_conf); tl != NULL;
175          tl = list_next(&tcps->tcps_listener_conf, tl)) {

```

```

175         if (tl->tl_port == lport) {
176             list_remove(&tcps->tcps_listener_conf, tl);
177             mutex_exit(&tcps->tcps_listener_conf_lock);
178             kmem_free(tl, sizeof (tcp_listener_t));
179             return (0);
180         }
181     }
182     mutex_exit(&tcps->tcps_listener_conf_lock);
183     return (ESRCH);
184 }
185
186 static int
187 tcp_set_buf_prop(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
188 const char *ifname, const void *pval, uint_t flags)
189 {
190     return (mod_set_buf_prop(stack->netstack_tcp->tcps_propinfo_tbl, stack,
191                             cr, pinfo, ifname, pval, flags));
192 }
193
194 static int
195 tcp_get_buf_prop(netstack_t *stack, mod_prop_info_t *pinfo, const char *ifname,
196 void *val, uint_t psize, uint_t flags)
197 {
198     return (mod_get_buf_prop(stack->netstack_tcp->tcps_propinfo_tbl, stack,
199                             pinfo, ifname, val, psize, flags));
200 }
201
202 /*
203  * Special checkers for smallest/largest anonymous port so they don't
204  * ever happen to be (largest < smallest).
205  */
206 /* ARGSUSED */
207 static int
208 tcp_smallest_anon_set(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
209 void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
210 const char *ifname, const void *pval, uint_t flags)
211 {
212     unsigned long new_value;
213     tcp_stack_t *tcps = stack->netstack_tcp;
214     tcp_stack_t *tcps = (tcp_stack_t *)cbarg;
215     int err;
216
217     if ((err = mod_uint32_value(pval, pinfo, flags, &new_value)) != 0)
218         return (err);
219     /* mod_uint32_value() + pinfo guarantees we're in TCP port range. */
220     if ((uint32_t)new_value > tcps->tcps_largest_anon_port)
221         return (ERANGE);
222     pinfo->prop_cur_uval = (uint32_t)new_value;
223     return (0);
224 }
225
226 static int
227 tcp_largest_anon_set(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
228 void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
229 const char *ifname, const void *pval, uint_t flags)
230 {
231     unsigned long new_value;
232     tcp_stack_t *tcps = stack->netstack_tcp;
233     tcp_stack_t *tcps = (tcp_stack_t *)cbarg;
234     int err;
235
236     if ((err = mod_uint32_value(pval, pinfo, flags, &new_value)) != 0)
237         return (err);
238     /* mod_uint32_value() + pinfo guarantees we're in TCP port range. */
239     if ((uint32_t)new_value < tcps->tcps_smallest_anon_port)

```

```

237         return (ERANGE);
238     pinfo->prop_cur_uval = (uint32_t)new_value;
239     return (0);
240 }

242 /*
243 * All of these are alterable, within the min/max values given, at run time.
244 *
245 * Note: All those tunables which do not start with "_" are Committed and
246 * therefore are public. See PSARC 2010/080.
247 */
248 mod_prop_info_t tcp_propinfo_tbl[] = {
249     /* tunable - 0 */
250     { "_time_wait_interval", MOD_PROTO_TCP,
251       mod_set_uint32, mod_get_uint32,
252       {1*SECONDS, 10*MINUTES, 1*MINUTES}, {1*MINUTES} },
253
254     { "_conn_req_max_q", MOD_PROTO_TCP,
255       mod_set_uint32, mod_get_uint32,
256       {1, UINT32_MAX, 128}, {128} },
257
258     { "_conn_req_max_q0", MOD_PROTO_TCP,
259       mod_set_uint32, mod_get_uint32,
260       {0, UINT32_MAX, 1024}, {1024} },
261
262     { "_conn_req_min", MOD_PROTO_TCP,
263       mod_set_uint32, mod_get_uint32,
264       {1, 1024, 1}, {1} },
265
266     { "_conn_grace_period", MOD_PROTO_TCP,
267       mod_set_uint32, mod_get_uint32,
268       {0*MS, 20*SECONDS, 0*MS}, {0*MS} },
269
270     { "_cwnd_max", MOD_PROTO_TCP,
271       mod_set_uint32, mod_get_uint32,
272       {128, ULP_MAX_BUF, 1024*1024}, {1024*1024} },
273
274     { "_debug", MOD_PROTO_TCP,
275       mod_set_uint32, mod_get_uint32,
276       {0, 10, 0}, {0} },
277
278     { "smallest_nonpriv_port", MOD_PROTO_TCP,
279       mod_set_uint32, mod_get_uint32,
280       {1024, (32*1024), 1024}, {1024} },
281
282     { "_ip_abort_cinterval", MOD_PROTO_TCP,
283       mod_set_uint32, mod_get_uint32,
284       {1*SECONDS, UINT32_MAX, 3*MINUTES}, {3*MINUTES} },
285
286     { "_ip_abort_linterval", MOD_PROTO_TCP,
287       mod_set_uint32, mod_get_uint32,
288       {1*SECONDS, UINT32_MAX, 3*MINUTES}, {3*MINUTES} },
289
290     /* tunable - 10 */
291     { "_ip_abort_interval", MOD_PROTO_TCP,
292       mod_set_uint32, mod_get_uint32,
293       {500*MS, UINT32_MAX, 5*MINUTES}, {5*MINUTES} },
294
295     { "_ip_notify_cinterval", MOD_PROTO_TCP,
296       mod_set_uint32, mod_get_uint32,
297       {1*SECONDS, UINT32_MAX, 10*SECONDS},
298       {10*SECONDS} },
299
300     { "_ip_notify_interval", MOD_PROTO_TCP,
301       mod_set_uint32, mod_get_uint32,

```

```

302       {500*MS, UINT32_MAX, 10*SECONDS}, {10*SECONDS} },
303
304     { "_ipv4_ttl", MOD_PROTO_TCP,
305       mod_set_uint32, mod_get_uint32,
306       {1, 255, 64}, {64} },
307
308     { "_keepalive_interval", MOD_PROTO_TCP,
309       mod_set_uint32, mod_get_uint32,
310       {10*SECONDS, 10*DAY, 2*HOURS}, {2*HOURS} },
311
312     { "_maxpsz_multiplier", MOD_PROTO_TCP,
313       mod_set_uint32, mod_get_uint32,
314       {0, 100, 10}, {10} },
315
316     { "_mss_def_ipv4", MOD_PROTO_TCP,
317       mod_set_uint32, mod_get_uint32,
318       {1, TCP_MSS_MAX_IPV4, 536}, {536} },
319
320     { "_mss_max_ipv4", MOD_PROTO_TCP,
321       mod_set_uint32, mod_get_uint32,
322       {1, TCP_MSS_MAX_IPV4, TCP_MSS_MAX_IPV4},
323       {TCP_MSS_MAX_IPV4} },
324
325     { "_mss_min", MOD_PROTO_TCP,
326       mod_set_uint32, mod_get_uint32,
327       {1, TCP_MSS_MAX, 108}, {108} },
328
329     { "_naglim_def", MOD_PROTO_TCP,
330       mod_set_uint32, mod_get_uint32,
331       {1, (64*1024)-1, (4*1024)-1}, {(4*1024)-1} },
332
333     /* tunable - 20 */
334     { "_rexmit_interval_initial", MOD_PROTO_TCP,
335       mod_set_uint32, mod_get_uint32,
336       {1*MS, 20*SECONDS, 1*SECONDS}, {1*SECONDS} },
337
338     { "_rexmit_interval_max", MOD_PROTO_TCP,
339       mod_set_uint32, mod_get_uint32,
340       {1*MS, 2*HOURS, 60*SECONDS}, {60*SECONDS} },
341
342     { "_rexmit_interval_min", MOD_PROTO_TCP,
343       mod_set_uint32, mod_get_uint32,
344       {1*MS, 2*HOURS, 400*MS}, {400*MS} },
345
346     { "_deferred_ack_interval", MOD_PROTO_TCP,
347       mod_set_uint32, mod_get_uint32,
348       {1*MS, 1*MINUTES, 100*MS}, {100*MS} },
349
350     { "_snd_lowat_fraction", MOD_PROTO_TCP,
351       mod_set_uint32, mod_get_uint32,
352       {0, 16, 10}, {10} },
353
354     { "_dupack_fast_retransmit", MOD_PROTO_TCP,
355       mod_set_uint32, mod_get_uint32,
356       {1, 10000, 3}, {3} },
357
358     { "_ignore_path_mtu", MOD_PROTO_TCP,
359       mod_set_boolean, mod_get_boolean,
360       {B_FALSE}, {B_FALSE} },
361
362     { "smallest_anon_port", MOD_PROTO_TCP,
363       tcp_smallest_anon_set, mod_get_uint32,
364       {1024, ULP_MAX_PORT, 32*1024}, {32*1024} },
365
366     { "largest_anon_port", MOD_PROTO_TCP,

```

```

367     tcp_largest_anon_set, mod_get_uint32,
368     {1024, ULP_MAX_PORT, ULP_MAX_PORT},
369     {ULP_MAX_PORT} },
370
371     { "send_buf", MOD_PROTO_TCP,
372       tcp_set_buf_prop, tcp_get_buf_prop,
373       {TCP_XMIT_LOWATER, ULP_MAX_BUF, TCP_XMIT_HIWATER},
374
375     { "send_maxbuf", MOD_PROTO_TCP,
376       mod_set_uint32, mod_get_uint32,
377       {TCP_XMIT_LOWATER, (1<<30), TCP_XMIT_HIWATER},
378       {TCP_XMIT_HIWATER} },
379
380     /* tunable - 30 */
381     { "_xmit_lowat", MOD_PROTO_TCP,
382       mod_set_uint32, mod_get_uint32,
383       {TCP_XMIT_LOWATER, ULP_MAX_BUF, TCP_XMIT_LOWATER},
384       {TCP_XMIT_LOWATER, (1<<30), TCP_XMIT_LOWATER},
385       {TCP_XMIT_LOWATER} },
386
387     { "recv_buf", MOD_PROTO_TCP,
388       tcp_set_buf_prop, tcp_get_buf_prop,
389       {TCP_RECV_LOWATER, ULP_MAX_BUF, TCP_RECV_HIWATER},
390
391     { "recv_maxbuf", MOD_PROTO_TCP,
392       mod_set_uint32, mod_get_uint32,
393       {TCP_RECV_LOWATER, (1<<30), TCP_RECV_HIWATER},
394       {TCP_RECV_HIWATER} },
395
396     { "_recv_hiwat_minmss", MOD_PROTO_TCP,
397       mod_set_uint32, mod_get_uint32,
398       {1, 65536, 4}, {4} },
399
400     /* Question: What default value should I set for tcp_strong_iss?
401    */
402     { "_strong_iss", MOD_PROTO_TCP,
403       mod_set_uint32, mod_get_uint32,
404       {0, 2, 1}, {1} },
405
406     { "_rtt_updates", MOD_PROTO_TCP,
407       mod_set_uint32, mod_get_uint32,
408       {0, 65536, 20}, {20} },
409
410     { "_wscale_always", MOD_PROTO_TCP,
411       mod_set_boolean, mod_get_boolean,
412       {B_TRUE}, {B_TRUE} },
413
414     { "_tstamp_always", MOD_PROTO_TCP,
415       mod_set_boolean, mod_get_boolean,
416       {B_FALSE}, {B_FALSE} },
417
418     { "_tstamp_if_wscale", MOD_PROTO_TCP,
419       mod_set_boolean, mod_get_boolean,
420       {B_TRUE}, {B_TRUE} },
421
422     /* tunable - 40 */

```

```

424     { "_rexmit_interval_extra", MOD_PROTO_TCP,
425       mod_set_uint32, mod_get_uint32,
426       {0*MS, 2*HOURS, 0*MS}, {0*MS} },
427
428     { "_deferred_acks_max", MOD_PROTO_TCP,
429       mod_set_uint32, mod_get_uint32,
430       {0, 16, 2}, {2} },
431
432     { "_slow_start_after_idle", MOD_PROTO_TCP,
433       mod_set_uint32, mod_get_uint32,
434       {0, 16384, 0}, {0} },
435
436     { "_slow_start_initial", MOD_PROTO_TCP,
437       mod_set_uint32, mod_get_uint32,
438       {0, 16, 0}, {0} },
439
440     { "sack", MOD_PROTO_TCP,
441       mod_set_uint32, mod_get_uint32,
442       {0, 2, 2}, {2} },
443
444     { "_ipv6_hoplimit", MOD_PROTO_TCP,
445       mod_set_uint32, mod_get_uint32,
446       {0, IPV6_MAX_HOPS, IPV6_DEFAULT_HOPS},
447       {IPV6_DEFAULT_HOPS} },
448
449     { "_mss_def_ipv6", MOD_PROTO_TCP,
450       mod_set_uint32, mod_get_uint32,
451       {1, TCP_MSS_MAX_IPV6, 1220}, {1220} },
452
453     { "_mss_max_ipv6", MOD_PROTO_TCP,
454       mod_set_uint32, mod_get_uint32,
455       {1, TCP_MSS_MAX_IPV6, TCP_MSS_MAX_IPV6},
456       {TCP_MSS_MAX_IPV6} },
457
458     { "_rev_src_routes", MOD_PROTO_TCP,
459       mod_set_boolean, mod_get_boolean,
460       {B_FALSE}, {B_FALSE} },
461
462     { "_local_dack_interval", MOD_PROTO_TCP,
463       mod_set_uint32, mod_get_uint32,
464       {10*MS, 500*MS, 50*MS}, {50*MS} },
465
466     /* tunable - 50 */
467     { "_local_dacks_max", MOD_PROTO_TCP,
468       mod_set_uint32, mod_get_uint32,
469       {0, 16, 8}, {8} },
470
471     { "ecn", MOD_PROTO_TCP,
472       mod_set_uint32, mod_get_uint32,
473       {0, 2, 1}, {1} },
474
475     { "_rst_sent_rate_enabled", MOD_PROTO_TCP,
476       mod_set_boolean, mod_get_boolean,
477       {B_TRUE}, {B_TRUE} },
478
479     { "_rst_sent_rate", MOD_PROTO_TCP,
480       mod_set_uint32, mod_get_uint32,
481       {0, UINT32_MAX, 40}, {40} },
482
483     { "_push_timer_interval", MOD_PROTO_TCP,
484       mod_set_uint32, mod_get_uint32,
485       {0, 100*MS, 50*MS}, {50*MS} },
486
487     { "_use_smss_as_mss_opt", MOD_PROTO_TCP,
488       mod_set_boolean, mod_get_boolean,
489       {B_FALSE}, {B_FALSE} },

```

```
491     { "_keepalive_abort_interval", MOD_PROTO_TCP,
492         mod_set_uint32, mod_get_uint32,
493         {0, UINT32_MAX, 8*MINUTES}, {8*MINUTES} },
494
495     /*  
496      * tcp_wroff_xtra is the extra space in front of TCP/IP header for link  
497      * layer header. It has to be a multiple of 8.  
498      */  
499     { "_wroff_xtra", MOD_PROTO_TCP,
500         mod_set_aligned, mod_get_uint32,
501         {0, 256, 32}, {32} },
502
503     { "_dev_flow_ctl", MOD_PROTO_TCP,
504         mod_set_boolean, mod_get_boolean,
505         {B_FALSE}, {B_FALSE} },
506
507     { "_reass_timeout", MOD_PROTO_TCP,
508         mod_set_uint32, mod_get_uint32,
509         {0, UINT32_MAX, 100*SECONDS}, {100*SECONDS} },
510
511     /* tunable - 60 */
512     { "extra_priv_ports", MOD_PROTO_TCP,
513         mod_set_extra_privports, mod_get_extra_privports,
514         {1, ULP_MAX_PORT, 0}, {0} },
515
516     { "_1948_phrase", MOD_PROTO_TCP,
517         tcp_set_1948phrase, NULL, {0}, {0} },
518
519     { "_listener_limit_conf",
520         NULL, tcp_listener_conf_get, {0}, {0} },
521
522     { "_listener_limit_conf_add",
523         NULL, tcp_listener_conf_add, {0}, {0} },
524
525     { "_listener_limit_conf_del",
526         NULL, tcp_listener_conf_del, {0}, {0} },
527
528     { "_iss_incr", MOD_PROTO_TCP,
529         mod_set_uint32, mod_get_uint32,
530         {1, ISS_INCR, ISS_INCR},
531         {ISS_INCR} },
532
533     { "?", MOD_PROTO_TCP, NULL, mod_get_allprop, {0}, {0} },
534
535     { NULL, 0, NULL, NULL, {0}, {0} }
536 };
```

unchanged portion omitted

```
new/usr/src/uts/common/inet/tcp_impl.h
```

```
*****  
30072 Mon Jul 29 18:04:24 2013
```

```
new/usr/src/uts/common/inet/tcp_impl.h
```

```
3942 inject sanity into ipadn tcp buffer size properties
```

```
3943 _snd_lowat_fraction tcp tunable has no effect
```

```
Reviewed by: Adam Leventhal <ahl@delphix.com>
```

```
Reviewed by: Peng Dai <peng.dai@delphix.com>
```

```
*****
```

```
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) 2005, 2010, Oracle and/or its affiliates. All rights reserved.  
23 * Copyright (c) 2011, Joyent Inc. All rights reserved.  
24 * Copyright (c) 2013, OmniTI Computer Consulting, Inc. All rights reserved.  
25 * Copyright (c) 2013 by Delphix. All rights reserved.  
26 */
```

```
28 #ifndef _INET_TCP_IMPL_H  
29 #define _INET_TCP_IMPL_H
```

```
31 /*  
32 * TCP implementation private declarations. These interfaces are  
33 * used to build the IP module and are not meant to be accessed  
34 * by any modules except IP itself. They are undocumented and are  
35 * subject to change without notice.  
36 */
```

```
38 #ifdef __cplusplus  
39 extern "C" {  
40 #endif
```

```
42 #ifdef _KERNEL
```

```
44 #include <sys/cpuvar.h>  
45 #include <sys/clock_impl.h> /* For LBOLT_FASTPATH{,64} */  
46 #include <inet/optcom.h>  
47 #include <inet/tcp.h>  
48 #include <inet/tunables.h>
```

```
50 #define TCP_MOD_ID      5105
```

```
52 extern struct qinit    tcp_sock_winit;  
53 extern struct qinit    tcp_winit;
```

```
55 extern sock_downcalls_t sock_tcp_downcalls;
```

```
57 /*  
58 * Note that by default, the _snd_lowat_fraction tunable controls the value of
```

```
1
```

```
new/usr/src/uts/common/inet/tcp_impl.h
```

```
59 * the transmit low water mark. TCP_XMIT_LOWATER (and thus the _xmit_lowat  
60 * property) is only used if the administrator has disabled _snd_lowat_fraction  
61 * by setting it to 0.  
62 */  
63 #define TCP_XMIT_LOWATER      4096  
64 #define TCP_XMIT_HIWATER     49152  
65 #define TCP_RECV_LOWATER     2048  
66 #define TCP_RECV_HIWATER    128000  
68 /*  
69 * Bind hash list size and has function. It has to be a power of 2 for  
70 * hashing.  
71 */  
72 #define TCP_BIND_FANOUT_SIZE   1024  
73 #define TCP_BIND_HASH(lport) (ntohs(lport) & (TCP_BIND_FANOUT_SIZE - 1))  
75 /*  
76 * This implementation follows the 4.3BSD interpretation of the urgent  
77 * pointer and not RFC 1122. Switching to RFC 1122 behavior would cause  
78 * incompatible changes in protocols like telnet and rlogin.  
79 */  
80 #define TCP_OLD_UPR_INTERPRETATION    1  
82 /* TCP option length */  
83 #define TCPOPT_NOP_LEN           1  
84 #define TCPOPT_MAXSEG_LEN        4  
85 #define TCPOPT_WS_LEN            3  
86 #define TCPOPT_REAL_WS_LEN       (TCPOPT_WS_LEN+1)  
87 #define TCPOPT_TSTAMP_LEN        10  
88 #define TCPOPT_REAL_TS_LEN       (TCPOPT_TSTAMP_LEN+2)  
89 #define TCPOPT_SACK_OK_LEN       2  
90 #define TCPOPT_REAL_SACK_OK_LEN (TCPOPT_SACK_OK_LEN+2)  
91 #define TCPOPT_REAL_SACK_LEN     4  
92 #define TCPOPT_MAX_SACK_LEN      36  
93 #define TCPOPT_HEADER_LEN         2  
95 /* Round up the value to the nearest mss. */  
96 #define MSS_ROUNDUP(value, mss) ((((value) - 1) / (mss) + 1) * (mss))  
98 /*  
99 * Was this tcp created via socket() interface?  
100 */  
101 #define TCP_IS_SOCKET(tcp) ((tcp)->tcp_issocket)  
103 /*  
104 * Is this tcp not attached to any upper client?  
105 */  
106 #define TCP_IS_DETACHED(tcp) ((tcp)->tcp_detached)  
108 /* TCP timers related data structures. Refer to tcp_timers.c. */  
109 typedef struct tcp_timer_s {  
110     conn_t *connp;  
111     void (*tcpt_proc)(void *);  
112     callout_id_t tcpt_tid;  
113 } tcp_timer_t;  
114 unchanged portion omitted
```

```
2
```

```
*****
13690 Mon Jul 29 18:04:25 2013
new/usr/src/uts/common/inet/tunables.c
3942 inject sanity into ipadm tcp buffer size properties
3943 _snd_lowat_fraction tcp tunable has no effect
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Peng Dai <peng.dai@delphix.com>
*****
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) 1991, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 1990 Mentor Inc.
24 * Copyright (c) 2013 by Delphix. All rights reserved.
25 */
26
27 #include <inet/tunables.h>
28 #include <sys/md5.h>
29 #include <inet/common.h>
30 #include <inet/ip.h>
31 #include <inet/ip6.h>
32 #include <netinet/icmp6.h>
33 #include <inet/ip_stack.h>
34 #include <inet/rawip_impl.h>
35 #include <inet/tcp_stack.h>
36 #include <inet/tcp_impl.h>
37 #include <inet/udp_impl.h>
38 #include <inet/sctp/sctp_stack.h>
39 #include <inet/sctp/sctp_impl.h>
40 #include <inet/tunables.h>
41
42 mod_prop_info_t *
43 mod_prop_lookup(mod_prop_info_t ptbl[], const char *prop_name, uint_t proto)
44 {
45     mod_prop_info_t *pinfo;
46
47     /*
48      * Walk the ptbl array looking for a property that has the requested
49      * name and protocol number. Note that we assume that all protocol
50      * tables are terminated by an entry with a NULL property name.
51      */
52     for (pinfo = ptbl; pinfo->mpi_name != NULL; pinfo++) {
53         if (strcmp(pinfo->mpi_name, prop_name) == 0 &&
54             pinfo->mpi_proto == proto)
55             return (pinfo);
56     }
57     return (NULL);
58 }
```

```
60 static int
61 prop_perm2const(mod_prop_info_t *pinfo)
62 {
63     if (pinfo->mpi_setf == NULL)
64         return (MOD_PROP_PERM_READ);
65     if (pinfo->mpi_getf == NULL)
66         return (MOD_PROP_PERM_WRITE);
67     return (MOD_PROP_PERM_RW);
68 }
69
70 /*
71  * Modifies the value of the property to default value or to the 'pval'
72  * specified by the user.
73  */
74 /* ARGSUSED */
75 int
76 mod_set_boolean(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
77 mod_set_boolean(void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
78 const char *ifname, const void* pval, uint_t flags)
79 {
80     char          *end;
81     unsigned long  new_value;
82
83     if (flags & MOD_PROP_DEFAULT) {
84         pinfo->prop_cur_bval = pinfo->prop_def_bval;
85         return (0);
86     }
87
88     if (ddi_strtoul(pval, &end, 10, &new_value) != 0 || *end != '\0')
89         return (EINVAL);
90     if (new_value != B_TRUE && new_value != B_FALSE)
91         return (EINVAL);
92     pinfo->prop_cur_bval = new_value;
93     return (0);
94 }
95 /*
96  * Retrieves property permission, default value, current value or possible
97  * values for those properties whose value type is boolean_t.
98 */
99 /* ARGSUSED */
100 int
101 mod_get_boolean(netstack_t *stack, mod_prop_info_t *pinfo, const char *ifname,
102 mod_get_boolean(void *cbarg, mod_prop_info_t *pinfo, const char *ifname,
103 void *pval, uint_t psize, uint_t flags)
104 {
105     boolean_t      get_def = (flags & MOD_PROP_DEFAULT);
106     boolean_t      get_perm = (flags & MOD_PROP_PERM);
107     boolean_t      get_range = (flags & MOD_PROP_POSSIBLE);
108     size_t         nbytes;
109
110     bzero(pval, psize);
111     if (get_perm)
112         nbytes = sprintf(pval, psize, "%u", prop_perm2const(pinfo));
113     else if (get_range)
114         nbytes = sprintf(pval, psize, "%u,%u", B_FALSE, B_TRUE);
115     else if (get_def)
116         nbytes = sprintf(pval, psize, "%u", pinfo->prop_def_bval);
117     else
118         nbytes = sprintf(pval, psize, "%u", pinfo->prop_cur_bval);
119     if (nbytes >= psize)
120         return (ENOBUFS);
121     return (0);
122 }
```

unchanged portion omitted

```

144 /*
145  * Modifies the value of the property to default value or to the 'pval'
146  * specified by the user.
147 */
148 /* ARGSUSED */
149 int
150 mod_set_uint32(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
151 mod_set_uint32(void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
152 {
153     const char *ifname, const void *pval, uint_t flags)
154 {
155     unsigned long new_value;
156     int err;
157
158     if ((err = mod_uint32_value(pval, pinfo, flags, &new_value)) != 0)
159         return (err);
160     pinfo->prop_cur_uval = (uint32_t)new_value;
161     return (0);
162 }
163 /* Rounds up the value to make it multiple of 8.
164 */
165 /* ARGSUSED */
166 int
167 mod_set_aligned(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
168 mod_set_aligned(void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
169 {
170     const char *ifname, const void* pval, uint_t flags)
171 {
172     int err;
173
174     if ((err = mod_set_uint32(stack, cr, pinfo, ifname, pval, flags)) != 0)
175     if ((err = mod_set_uint32(cbarg, cr, pinfo, ifname, pval, flags)) != 0)
176         return (err);
177
178     /* if required, align the value to multiple of 8 */
179     if (pinfo->prop_cur_uval & 0x7) {
180         pinfo->prop_cur_uval &= ~0x7;
181         pinfo->prop_cur_uval += 0x8;
182     }
183
184 /*
185  * Retrieves property permission, default value, current value or possible
186  * values for those properties whose value type is uint32_t.
187 */
188 /* ARGSUSED */
189 int
190 mod_get_uint32(netstack_t *stack, mod_prop_info_t *pinfo, const char *ifname,
191 mod_get_uint32(void *cbarg, mod_prop_info_t *pinfo, const char *ifname,
192 {
193     boolean_t get_def = (flags & MOD_PROP_DEFAULT);
194     boolean_t get_perm = (flags & MOD_PROP_PERM);
195     boolean_t get_range = (flags & MOD_PROP_POSSIBLE);
196     size_t nbytes;
197
198     bzero(pval, psizes);
199     if (get_perm)
200         nbytes = snprintf(pval, psizes, "%u", prop_perm2const(pinfo));
201     else if (get_range)
202         nbytes = snprintf(pval, psizes, "%u-%u",
203                           pinfo->prop_min_uval, pinfo->prop_max_uval);
204     else if (get_def)

```

```

205             nbytes = sprintf(pval, psizes, "%u", pinfo->prop_def_uval);
206         else
207             nbytes = sprintf(pval, psizes, "%u", pinfo->prop_cur_uval);
208         if (nbytes >= psizes)
209             return (ENOBUFS);
210         return (0);
211     }
212
213 /*
214  * The range of the buffer size properties has a static lower bound configured
215  * in the property info structure of the property itself, and a dynamic upper
216  * bound. The upper bound is the current value of the "max_buf" property
217  * in the appropriate protocol property table.
218 */
219 static void
220 mod_get_buf_prop_range(mod_prop_info_t ptbl[], mod_prop_info_t *pinfo,
221                         uint32_t *min, uint32_t *max)
222 {
223     mod_prop_info_t *maxbuf_pinfo = mod_prop_lookup(ptbl, "max_buf",
224                                                     pinfo->mpi_proto);
225
226     *min = pinfo->prop_min_uval;
227     *max = maxbuf_pinfo->prop_cur_uval;
228 }
229
230 /*
231  * Modifies the value of the buffer size property to its default value or to
232  * the value specified by the user. This is similar to mod_set_uint32() except
233  * that the value has a dynamically bounded range (see mod_get_buf_prop_range())
234  * for details).
235 */
236 /* ARGSUSED */
237 int
238 mod_set_buf_prop(mod_prop_info_t ptbl[], netstack_t *stack, cred_t *cr,
239                   mod_prop_info_t *pinfo, const char *ifname, const void *pval, uint_t flags)
240 {
241     unsigned long new_value;
242     char *end;
243     uint32_t min, max;
244
245     if (flags & MOD_PROP_DEFAULT) {
246         pinfo->prop_cur_uval = pinfo->prop_def_uval;
247         return (0);
248     }
249
250     if (ddi_strtoul(pval, &end, 10, &new_value) != 0 || *end != '\0')
251         return (EINVAL);
252
253     mod_get_buf_prop_range(ptbl, pinfo, &min, &max);
254     if (new_value < min || new_value > max)
255         return (ERANGE);
256
257     pinfo->prop_cur_uval = new_value;
258     return (0);
259 }
260
261 /*
262  * Retrieves property permissions, default value, current value, or possible
263  * values for buffer size properties. While these properties have integer
264  * values, they have a dynamic range (see mod_get_buf_prop_range() for
265  * details). As such, they need to be handled differently.
266 */
267 int
268 mod_get_buf_prop(mod_prop_info_t ptbl[], netstack_t *stack,
269                   mod_prop_info_t *pinfo, const char *ifname, void *pval, uint_t psizes,
270                   uint_t flags)

```

```

271 {
272     size_t nbytes;
273     uint32_t min, max;
274
275     if (flags & MOD_PROP_POSSIBLE) {
276         mod_get_buf_prop_range(ptbl, pinfo, &min, &max);
277         nbytes = snprintf(pval, psize, "%u-%u", min, max);
278         return (nbytes < psize ? 0 : ENOBUFS);
279     }
280     return (mod_get_uint32(stack, pinfo, ifname, pval, psize, flags));
281 }
282 */
283 * Implements /sbin/ndd -get /dev/ip ?, for all the modules. Needed for
284 * backward compatibility with /sbin/ndd.
285 */
286 /* ARGSUSED */
287 int
288 mod_get_allprop(netstack_t *stack, mod_prop_info_t *pinfo, const char *ifname,
289 mod_get_allprop(void *cbarg, mod_prop_info_t *pinfo, const char *ifname,
290 void *val, uint_t psize, uint_t flags)
291 {
292     char *pval = val;
293     mod_prop_info_t *ptbl, *prop;
294     ip_stack_t *ipst;
295     tcp_stack_t *tcps;
296     sctp_stack_t *sctps;
297     udp_stack_t *us;
298     icmp_stack_t *is;
299     uint_t size;
300     size_t nbytes = 0, tbytes = 0;
301
302     bzero(pval, psize);
303     size = psize;
304
305     switch (pinfo->mpi_proto) {
306     case MOD_PROTO_IP:
307     case MOD_PROTO_IPV4:
308     case MOD_PROTO_IPV6:
309         ptbl = stack->netstack_ip->ips_propinfo_tbl;
310         ipst = (ip_stack_t *)cbarg;
311         ptbl = ipst->ips_propinfo_tbl;
312         break;
313     case MOD_PROTO_RAWIP:
314         ptbl = stack->netstack_icmp->is_propinfo_tbl;
315         is = (icmp_stack_t *)cbarg;
316         ptbl = is->is_propinfo_tbl;
317         break;
318     case MOD_PROTO_TCP:
319         ptbl = stack->netstack_tcp->tcps_propinfo_tbl;
320         tcps = (tcp_stack_t *)cbarg;
321         ptbl = tcps->tcps_propinfo_tbl;
322         break;
323     case MOD_PROTO_UDP:
324         ptbl = stack->netstack_udp->us_propinfo_tbl;
325         us = (udp_stack_t *)cbarg;
326         ptbl = us->us_propinfo_tbl;
327         break;
328     case MOD_PROTO_SCTP:
329         ptbl = stack->netstack_sctp->sctps_propinfo_tbl;
330         sctps = (sctp_stack_t *)cbarg;
331         ptbl = sctps->sctps_propinfo_tbl;
332         break;
333     default:
334         return (EINVAL);
335     }

```

```

322     for (prop = ptbl; prop->mpi_name != NULL; prop++) {
323         if (prop->mpi_name[0] == '\0' ||
324             strcmp(prop->mpi_name, "?") == 0) {
325             continue;
326         }
327         nbytes = sprintf(pval, size, "%s %d %d", prop->mpi_name,
328                           prop->mpi_proto, prop_perm2const(prop));
329         size -= nbytes + 1;
330         pval += nbytes + 1;
331         tbytes += nbytes + 1;
332         if (tbytes >= psize) {
333             /* Buffer overflow, stop copying information */
334             return (ENOBUFS);
335         }
336     }
337     return (0);
338 }
339 */
340 */
341 * Hold a lock while changing *_epriv_ports to prevent multiple
342 * threads from changing it at the same time.
343 */
344 /* ARGSUSED */
345 int
346 mod_set_extra_privports(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
347 mod_set_extra_privports(void *cbarg, cred_t *cr, mod_prop_info_t *pinfo,
348 const char *ifname, const void* val, uint_t flags)
349 {
350     uint_t proto = pinfo->mpi_proto;
351     tcp_stack_t *tcps;
352     sctp_stack_t *sctps;
353     udp_stack_t *us;
354     unsigned long new_value;
355     char *end;
356     kmutex_t *lock;
357     uint_t i, nports;
358     in_port_t *ports;
359     boolean_t def = (flags & MOD_PROP_DEFAULT);
360     const char *pval = val;
361
362     if (!def) {
363         if (ddi_strtoul(pval, &end, 10, &new_value) != 0 ||
364             *end != '\0') {
365             return (EINVAL);
366         }
367         if (new_value < pinfo->prop_min_uval ||
368             new_value > pinfo->prop_max_uval) {
369             return (ERANGE);
370         }
371     }
372
373     switch (proto) {
374     case MOD_PROTO_TCP:
375         tcps = stack->netstack_tcp;
376         tcps = (tcp_stack_t *)cbarg;
377         lock = &tcps->tcps_epriv_port_lock;
378         ports = tcps->tcps_g_epriv_ports;
379         nports = tcps->tcps_g_num_epriv_ports;
380         break;
381     case MOD_PROTO_UDP:
382         us = stack->netstack_udp;
383         us = (udp_stack_t *)cbarg;
384         lock = &us->us_epriv_port_lock;
385         ports = us->us_epriv_ports;

```

```

384         nports = us->us_num_epriv_ports;
385         break;
386     case MOD_PROTO_SCTP:
387         sctps = stack->netstack_sctp;
388         sctps = (sctp_stack_t *)cbarg;
389         lock = &sctps->sctps_epriv_port_lock;
390         ports = sctps->sctps_g_epriv_ports;
391         nports = sctps->sctps_g_num_epriv_ports;
392         break;
393     default:
394         return (ENOTSUP);
395     }
396
397     mutex_enter(lock);
398
399     /* if MOD_PROP_DEFAULT is set then reset the ports list to default */
400     if (def) {
401         for (i = 0; i < nports; i++)
402             ports[i] = 0;
403         ports[0] = ULP_DEF_EPRIV_PORT1;
404         ports[1] = ULP_DEF_EPRIV_PORT2;
405         mutex_exit(lock);
406         return (0);
407     }
408
409     /* Check if the value is already in the list */
410     for (i = 0; i < nports; i++) {
411         if (new_value == ports[i])
412             break;
413     }
414
415     if (flags & MOD_PROP_REMOVE) {
416         if (i == nports)
417             mutex_exit(lock);
418         return (ESRCH);
419     }
420     /* Clear the value */
421     ports[i] = 0;
422 } else if (flags & MOD_PROP_APPEND) {
423     if (i != nports)
424         mutex_exit(lock);
425     return (EEXIST);
426 }
427
428     /* Find an empty slot */
429     for (i = 0; i < nports; i++) {
430         if (ports[i] == 0)
431             break;
432     }
433     if (i == nports) {
434         mutex_exit(lock);
435         return (EOVERFLOW);
436     }
437     /* Set the new value */
438     ports[i] = (in_port_t)new_value;
439 } else {
440     /*
441      * If the user used 'assignment' modifier.
442      * For eg:
443      *      # ipadm set-prop -p extra_priv_ports=3001 tcp
444      *
445      * We clear all the ports and then just add 3001.
446      */
447     ASSERT(flags == MOD_PROP_ACTIVE);
448     for (i = 0; i < nports; i++)
449         ports[i] = 0;

```

```

449                     ports[0] = (in_port_t)new_value;
450     }
451
452     mutex_exit(lock);
453     return (0);
454 }
455 /*
456  * Note: No locks are held when inspecting *_epriv_ports
457  * but instead the code relies on:
458  * - the fact that the address of the array and its size never changes
459  * - the atomic assignment of the elements of the array
460 */
461 /* ARGSUSED */
462 int
463 mod_get_extra_privports(netstack_t *stack, mod_prop_info_t *pinfo,
464                         const char *ifname, void *val, uint_t psize, uint_t flags)
465 mod_get_extra_privports(void *cbarg, mod_prop_info_t *pinfo, const char *ifname,
466                         void *val, uint_t psize, uint_t flags)
467 {
468     uint_t          proto = pinfo->mpi_proto;
469     tcp_stack_t    *tcp;
470     sctp_stack_t   *sctps;
471     udp_stack_t    *us;
472     uint_t          i, nports, size;
473     in_port_t       *ports;
474     char            *pval = val;
475     size_t          nbytes = 0, tbytes = 0;
476     boolean_t       get_def = (flags & MOD_PROP_DEFAULT);
477     boolean_t       get_perm = (flags & MOD_PROP_PERM);
478     boolean_t       get_range = (flags & MOD_PROP_POSSIBLE);
479     bzero(pval, psize);
480     size = psize;
481
482     if (get_def) {
483         tbytes = sprintf(pval, psize, "%u,%u", ULP_DEF_EPRIV_PORT1,
484                           ULP_DEF_EPRIV_PORT2);
485         goto ret;
486     } else if (get_perm) {
487         tbytes = sprintf(pval, psize, "%u", MOD_PROP_PERM_RW);
488         goto ret;
489     }
490
491     switch (proto) {
492     case MOD_PROTO_TCP:
493         tcp = stack->netstack_tcp;
494         cbarg = (tcp_stack_t *)cbarg;
495         ports = tcp->tcps_g_epriv_ports;
496         nports = tcp->tcps_g_num_epriv_ports;
497         break;
498     case MOD_PROTO_UDP:
499         us = stack->netstack_udp;
500         cbarg = (udp_stack_t *)cbarg;
501         ports = us->us_epriv_ports;
502         nports = us->us_num_epriv_ports;
503         break;
504     case MOD_PROTO_SCTP:
505         sctps = stack->netstack_sctp;
506         cbarg = (sctp_stack_t *)cbarg;
507         ports = sctps->sctps_g_epriv_ports;
508         nports = sctps->sctps_g_num_epriv_ports;
509         break;
510     default:
511         return (ENOTSUP);
512     }

```

```
511     if (get_range) {
512         tbytes = snprintf(pval, psize, "%u-%u", pinfo->prop_min_uval,
513                           pinfo->prop_max_uval);
514         goto ret;
515     }
516
517     for (i = 0; i < nports; i++) {
518         if (ports[i] != 0) {
519             if (psize == size)
520                 nbytes = sprintf(pval, size, "%u", ports[i]);
521             else
522                 nbytes = sprintf(pval, size, ",%u", ports[i]);
523             size -= nbytes;
524             pval += nbytes;
525             tbytes += nbytes;
526             if (tbytes >= psizes)
527                 return (ENOBUFS);
528         }
529     }
530     return (0);
531 ret:
532     if (tbytes >= psizes)
533         return (ENOBUFS);
534     return (0);
535 }
```

unchanged portion omitted

```
*****
6596 Mon Jul 29 18:04:27 2013
new/usr/src/uts/common/inet/tunables.h
3942 inject sanity into ipadmin tcp buffer size properties
3943 _snd_lowat_fraction tcp tunable has no effect
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Peng Dai <peng.dai@delphix.com>
*****
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) 1991, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 1990 Mentor Inc.
24 * Copyright (c) 2013 by Delphix. All rights reserved.
25 */
26
27 #ifndef _INET_TUNABLES_H
28 #define _INET_TUNABLES_H
29
30 #include <sys/types.h>
31 #include <net/if.h>
32 #ifdef _KERNEL
33 #include <sys/netstack.h>
34 #endif
35
36 #ifdef __cplusplus
37 extern "C" {
38 #endif
39
40 #define MAXPROPNAMELEN 64
41
42 /*
43 * The 'mod_ioc_prop_s' datastructure is used as an IOCTL argument for
44 * SIOCSETPROP and SIOCGETPROP ioctls. This datastructure identifies the
45 * protocol ('mpr_proto') property ('mpr_name'), which needs to be modified
46 * or retrieved ('mpr_valsize' and 'mpr_val'). If the property applies to an
47 * interface then 'mpr_ifname' contains the name of the interface.
48 */
49 typedef struct mod_ioc_prop_s {
50     uint_t          mpr_version;
51     uint_t          mpr_flags;           /* see below */
52     /* name of the interface (ill) for which property will be applied */
53     char            mpr_ifname[LIFNAMSIZ];
54     uint_t          mpr_proto;          /* see below */
55     char            mpr_name[MAXPROPNAMELEN]; /* property name */
56     uint_t          mpr_valsize;        /* size of mpr_val */
57     char            mpr_val[1];
58 } mod_ioc_prop_t;
```

```
60 #define MOD_PROP_VERSION      1
61 /* permission flags for properties */
62 #define MOD_PROP_PERM_READ    0x1
63 #define MOD_PROP_PERM_WRITE   0x2
64 #define MOD_PROP_PERM_RW      (MOD_PROP_PERM_READ|MOD_PROP_PERM_WRITE)
65
66 /* mpr_flags values */
67 #define MOD_PROP_ACTIVE       0x01 /* current value of the property */
68 #define MOD_PROP_DEFAULT      0x02 /* default value of the property */
69 #define MOD_PROP_POSSIBLE     0x04 /* possible values for the property */
70 #define MOD_PROP_PERM         0x08 /* read/write permission for property */
71 #define MOD_PROP_APPEND       0x10 /* append to multi-valued property */
72 #define MOD_PROP_REMOVE       0x20 /* remove from multi-valued property */
73
74 /* mpr_proto values */
75 #define MOD_PROTO_NONE        0x00 /* property is applicable to IPV4 */
76 #define MOD_PROTO_IPV4        0x01 /* property is applicable to IPV6 */
77 #define MOD_PROTO_IPV6        0x02 /* property is applicable to ICMP */
78 #define MOD_PROTO_RAWIP       0x04 /* property is applicable to TCP */
79 #define MOD_PROTO_TCP         0x08 /* property is applicable to UDP */
80 #define MOD_PROTO_UDP         0x10 /* property is applicable to SCTP */
81 #define MOD_PROTO_SCTP        0x20 /* property is applicable to SCTP */
82
83 /* property is applicable to both IPV4[6] */
84 #define MOD_PROTO_IP          (MOD_PROTO_IPV4|MOD_PROTO_IPV6)
85
86 #ifdef _KERNEL
87
88 typedef struct mod_prop_info_s mod_prop_info_t;
89
90 /* set/get property callback functions */
91 typedef int    mod_prop_set_t(netstack_t *, cred_t *, mod_prop_info_t *,
92                               const char *, const void *, uint_t);
93
94 typedef int    mod_prop_get_t(netstack_t *, mod_prop_info_t *, const char *,
95                               void *, uint_t, uint_t);
96
97 typedef struct mod_propval_uint32_s {
98     uint32_t      mod_propval_umin;
99     uint32_t      mod_propval_umax;
100    uint32_t      mod_propval_ucur;
101 } mod_propval_uint32_t;
102
103
104 /* shortcuts to access current/default values */
105 #define prop_min_uval u.mpi_uval.mod_propval_umin
106 #define prop_max_uval u.mpi_uval.mod_propval_umax
107 #define prop_cur_uval u.mpi_uval.mod_propval_ucur
108 #define prop_bval      u.mpi_bval
109 #define prop_def_uval u_def.mpi_def_uval
110 #define prop_def_bval u_def.mpi_def_bval
111
112
113 #define MS             1L
114 #define SECONDS        (1000 * MS)
115 #define MINUTES        (60 * SECONDS)
116 #define HOURS          (60 * MINUTES)
117 #define DAYS           (24 * HOURS)
118
119
120 #define MB             (1024 * 1024)
121
122
123 /* Largest TCP/UDP/SCTP port number */
124 #define ULP_MAX_PORT   (64 * 1024 - 1)
```

```
149 /* extra privilege ports for upper layer protocols, tcp, sctp and udp */
150 #define ULP_DEF_EPRIV_PORT1    2049
151 #define ULP_DEF_EPRIV_PORT2    4045
153 #define ULP_MAX_BUF      (1<<30) /* Largest possible send/receive buffer */
155 /* generic function to set/get global module properties */
156 extern mod_prop_setf_t  mod_set_boolean, mod_set_uint32,
157                      mod_set_aligned, mod_set_extra_privports;
159 extern mod_prop_getf_t  mod_get_boolean, mod_get_uint32,
160                      mod_get_allprop, mod_get_extra_privports;
162 extern int               mod_uint32_value(const void *, mod_prop_info_t *,
163     uint_t, unsigned long *);
164 extern mod_prop_info_t *mod_prop_lookup(mod_prop_info_t[], const char *,
165     uint_t);
166 extern int               mod_set_buf_prop(mod_prop_info_t[], netstack_t *,
167     cred_t *cr, mod_prop_info_t *, const char *, const void *, uint_t);
168 extern int               mod_get_buf_prop(mod_prop_info_t[], netstack_t *,
169     mod_prop_info_t *, const char *, void *, uint_t, uint_t);
159 extern int mod_uint32_value(const void *, mod_prop_info_t *, uint_t,
160     unsigned long );
171 #endif /* _KERNEL */
173 /*
174  * End-system model definitions that include the weak/strong end-system
175  * definitions in RFC 1122, Section 3.3.4.5. IP_WEAK_ES and IP_STRONG_ES
176  * conform to the corresponding RFC 1122 definitions. The IP_SRC_PRI_ES
177  * hostmodel is similar to IP_WEAK_ES with one additional enhancement: for
178  * a packet with source S2, destination D2, the route selection algorithm
179  * will first attempt to find a route for the destination that goes out
180  * through an interface where S2 is configured and marked UP. If such
181  * a route cannot be found, then the best-matching route for D2 will be
182  * selected, ignoring any mismatches between S2 and the interface addresses
183  * on the outgoing interface implied by the route.
184 */
185 typedef enum {
186     IP_WEAK_ES = 0,
187     IP_SRC_PRI_ES,
188     IP_STRONG_ES,
189     IP_MAXVAL_ES
190 } ip_hostmodel_t;


---

unchanged portion omitted
```

```
new/usr/src/uts/common/inet/udp/udp_tunables.c
```

```
*****
5002 Mon Jul 29 18:04:28 2013
new/usr/src/uts/common/inet/udp/udp_tunables.c
3942 inject sanity into ipadn tcp buffer size properties
3943 _snd_lowat_fraction tcp tunable has no effect
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Peng Dai <peng.dai@delphix.com>
*****
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) 1991, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright 2013 Nexenta Systems, Inc. All rights reserved.
24 * Copyright (c) 2013 by Delphix. All rights reserved.
25 */
26 /* Copyright (c) 1990 Mentor Inc. */

28 #include <inet/ip.h>
29 #include <inet/ip6.h>
30 #include <inet/udp_impl.h>
31 #include <sys/sunddi.h>

33 static int
34 udp_set_buf_prop(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
35     const char *ifname, const void *pval, uint_t flags)
36 {
37     return (mod_set_buf_prop(stack->netstack_udp->us_propinfo_tbl, stack,
38         cr, pinfo, ifname, pval, flags));
39 }

41 static int
42 udp_get_buf_prop(netstack_t *stack, mod_prop_info_t *pinfo, const char *ifname,
43     void *val, uint_t psize, uint_t flags)
44 {
45     return (mod_get_buf_prop(stack->netstack_udp->us_propinfo_tbl, stack,
46         pinfo, ifname, val, psize, flags));
47 }

49 /*
50 * Special checkers for smallest/largest anonymous port so they don't
51 * ever happen to be (largest < smallest).
52 */
53 /* ARGSUSED */
54 static int
55 udp_smallest_anon_set(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
56     const char *ifname, const void *pval, uint_t flags)
57 {
```

```
1
```

```
new/usr/src/uts/common/inet/udp/udp_tunables.c
```

```
58     unsigned long new_value;
59     udp_stack_t *us = stack->netstack_udp;
60     udp_stack_t *us = (udp_stack_t *)cbarg;
61     int err;

62     if ((err = mod_uint32_value(pval, pinfo, flags, &new_value)) != 0)
63         return (err);
64     /* mod_uint32_value() + pinfo guarantees we're in UDP port range. */
65     if (new_value > us->us_largest_anon_port)
66         return (ERANGE);
67     pinfo->prop_cur_uval = (uint32_t)new_value;
68 }
69 }

71 /* ARGSUSED */
72 static int
73 udp_largest_anon_set(netstack_t *stack, cred_t *cr, mod_prop_info_t *pinfo,
74     const char *ifname, const void *pval, uint_t flags)
75 {
76     unsigned long new_value;
77     udp_stack_t *us = stack->netstack_udp;
78     udp_stack_t *us = (udp_stack_t *)cbarg;
79     int err;

80     if ((err = mod_uint32_value(pval, pinfo, flags, &new_value)) != 0)
81         return (err);
82     /* mod_uint32_value() + pinfo guarantees we're in UDP port range. */
83     if (new_value < us->us_smallest_anon_port)
84         return (ERANGE);
85     pinfo->prop_cur_uval = (uint32_t)new_value;
86     return (0);
87 }

89 /*
90 * All of these are alterable, within the min/max values given, at run time.
91 *
92 * Note: All those tunables which do not start with "_" are Committed and
93 * therefore are public. See PSARC 2010/080.
94 */
95 mod_prop_info_t udp_propinfo_tbl[] = {
96     /* tunable - 0 */
97     { "_wroff_extra", MOD_PROTO_UDP,
98         mod_set_uint32, mod_get_uint32,
99         {0, 256, 32}, {32} },
100    { "_ipv4_ttl", MOD_PROTO_UDP,
101        mod_set_uint32, mod_get_uint32,
102        {1, 255, 255}, {255} },
103    { "_ipv6_hoplimit", MOD_PROTO_UDP,
104        mod_set_uint32, mod_get_uint32,
105        {0, IPV6_MAX_HOPS, IPV6_DEFAULT_HOPS}, {IPV6_DEFAULT_HOPS} },
106    { "smallest_nonpriv_port", MOD_PROTO_UDP,
107        mod_set_uint32, mod_get_uint32,
108        {1024, (32 * 1024), 1024}, {1024} },
109    { "_do_checksum", MOD_PROTO_UDP,
110        mod_set_boolean, mod_get_boolean,
111        {B_TRUE}, {B_TRUE} },
112    { "smallest_anon_port", MOD_PROTO_UDP,
113        udp_smallest_anon_set, mod_get_uint32,
114        {1024, ULP_MAX_PORT, (32 * 1024)}, {(32 * 1024)} },
115 }
```

```
2
```

```
121     { "largest_anon_port", MOD_PROTO_UDP,
122       udp_largest_anon_set, mod_get_uint32,
123       {1024, ULP_MAX_PORT, ULP_MAX_PORT}, {ULP_MAX_PORT} },
124
125     { "send_buf", MOD_PROTO_UDP,
126       udp_set_buf_prop, udp_get_buf_prop,
127       {UDP_XMIT_LOWATER, ULP_MAX_BUF, UDP_XMIT_HIWATER},
128     { "send_maxbuf", MOD_PROTO_UDP,
129       mod_set_uint32, mod_get_uint32,
130       {UDP_XMIT_LOWATER, (1<<30), UDP_XMIT_HIWATER},
131       {UDP_XMIT_HIWATER} },
132
133     { "_xmit_lowat", MOD_PROTO_UDP,
134       mod_set_uint32, mod_get_uint32,
135       {0, ULP_MAX_BUF, UDP_XMIT_LOWATER},
136       {0, (1<<30), UDP_XMIT_LOWATER},
137       {UDP_XMIT_LOWATER} },
138
139     { "recv_buf", MOD_PROTO_UDP,
140       udp_set_buf_prop, udp_get_buf_prop,
141       {UDP_RECV_LOWATER, ULP_MAX_BUF, UDP_RECV_HIWATER},
142     { "recv_maxbuf", MOD_PROTO_UDP,
143       mod_set_uint32, mod_get_uint32,
144       {UDP_RECV_LOWATER, (1<<30), UDP_RECV_HIWATER},
145       {UDP_RECV_HIWATER} },
146
147     /* tunable - 10 */
148     { "max_buf", MOD_PROTO_UDP,
149       "_max_buf", MOD_PROTO_UDP,
150       mod_set_uint32, mod_get_uint32,
151       {65536, ULP_MAX_BUF, 2*1024*1024}, {2*1024*1024} },
152       {65536, (1<<30), 2*1024*1024}, {2*1024*1024} },
153
154     { "_pmtu_discovery", MOD_PROTO_UDP,
155       mod_set_boolean, mod_get_boolean,
156       {B_FALSE}, {B_FALSE} },
157
158     { "_sendto_ignerr", MOD_PROTO_UDP,
159       mod_set_boolean, mod_get_boolean,
160       {B_FALSE}, {B_FALSE} },
161
162     { "extra_priv_ports", MOD_PROTO_UDP,
163       mod_set_extra_privports, mod_get_extra_privports,
164       {1, ULP_MAX_PORT, 0}, {0} },
165
166     { "?", MOD_PROTO_UDP, NULL, mod_get_allprop, {0}, {0} },
167
168     { NULL, 0, NULL, NULL, {0}, {0} }
169 },
170 }
```

unchanged portion omitted