```
*********************************************************
  238265 Tue Jan 15 10:29:45 2019
new/usr/src/uts/common/os/sunmdi.c
10094 i_mdi_client_free() doesn't need to check for a NULL cdip
*********************************************************
   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) 2000, 2010, Oracle and/or its affiliates. All rights reserved.
  23  * Copyright (c) 2014 Nexenta Systems Inc. All rights reserved.
  24  * Copyright (c) 2018, Joyent, Inc.
  25  */

  27 /*
  28  * Multipath driver interface (MDI) implementation; see mdi_impldefs.h for a
  29  * more detailed discussion of the overall mpxio architecture.
  30  *
  31  * Default locking order:
  32  *
  33  * _NOTE(LOCK_ORDER(mdi_mutex, mdi_vhci:vh_phci_mutex));
  34  * _NOTE(LOCK_ORDER(mdi_mutex, mdi_vhci:vh_client_mutex));
  35  * _NOTE(LOCK_ORDER(mdi_vhci:vh_phci_mutex, mdi_phci::ph_mutex);
  36  * _NOTE(LOCK_ORDER(mdi_vhci:vh_client_mutex, mdi_client::ct_mutex);
  37  * _NOTE(LOCK_ORDER(mdi_phci::ph_mutex mdi_pathinfo::pi_mutex))
  38  * _NOTE(LOCK_ORDER(mdi_phci::ph_mutex mdi_client::ct_mutex))
  39  * _NOTE(LOCK_ORDER(mdi_client::ct_mutex mdi_pathinfo::pi_mutex))
  40  */

  42 #include <sys/note.h>
  43 #include <sys/types.h>
  44 #include <sys/varargs.h>
  45 #include <sys/param.h>
  46 #include <sys/errno.h>
  47 #include <sys/uio.h>
  48 #include <sys/buf.h>
  49 #include <sys/modctl.h>
  50 #include <sys/open.h>
  51 #include <sys/kmem.h>
  52 #include <sys/poll.h>
  53 #include <sys/conf.h>
  54 #include <sys/bootconf.h>
  55 #include <sys/cmn_err.h>
  56 #include <sys/stat.h>
  57 #include <sys/ddi.h>
  58 #include <sys/sunddi.h>
  59 #include <sys/ddipropdefs.h>
  60 #include <sys/sunndi.h>
  61 #include <sys/ndi_impldefs.h>
```

```
  62 #include <sys/promif.h>
  63 #include <sys/sunmdi.h>
  64 #include <sys/mdi_impldefs.h>
  65 #include <sys/taskq.h>
  66 #include <sys/epm.h>
  67 #include <sys/sunpm.h>
  68 #include <sys/modhash.h>
  69 #include <sys/disp.h>
  70 #include <sys/autoconf.h>
  71 #include <sys/sysmacros.h>

  73 #ifdef  DEBUG
  74 #include <sys/debug.h>
  75 int     mdi_debug = 1;
  76 int     mdi_debug_logonly = 0;
  77 #define MDI_DEBUG(dbglevel, pargs) if (mdi_debug >= (dbglevel)) i_mdi_log pargs
  78 #define MDI_WARN         CE_WARN, __func__
  79 #define MDI_NOTE         CE_NOTE, __func__
  80 #define MDI_CONT         CE_CONT, __func__
  81 static void i_mdi_log(int, const char *, dev_info_t *, const char *, ...);
  82 #else   /* !DEBUG */
  83 #define MDI_DEBUG(dbglevel, pargs)
  84 #endif  /* DEBUG */
  85 int     mdi_debug_consoleonly = 0;
  86 int     mdi_delay = 3;

  88 extern pri_t    minclsyspri;
  89 extern int      modrootloaded;

  91 /*
  92  * Global mutex:
  93  * Protects vHCI list and structure members.
  94  */
  95 kmutex_t        mdi_mutex;

  97 /*
  98  * Registered vHCI class driver lists
  99  */
 100 int             mdi_vhci_count;
 101 mdi_vhci_t      *mdi_vhci_head;
 102 mdi_vhci_t      *mdi_vhci_tail;

 104 /*
 105  * Client Hash Table size
 106  */
 107 static int      mdi_client_table_size = CLIENT_HASH_TABLE_SIZE;

 109 /*
 110  * taskq interface definitions
 111  */
 112 #define MDI_TASKQ_N_THREADS     8
 113 #define MDI_TASKQ_PRI           minclsyspri
 114 #define MDI_TASKQ_MINALLOC      (4*mdi_taskq_n_threads)
 115 #define MDI_TASKQ_MAXALLOC      (500*mdi_taskq_n_threads)

 117 taskq_t                         *mdi_taskq;
 118 static uint_t                   mdi_taskq_n_threads = MDI_TASKQ_N_THREADS;

 120 #define TICKS_PER_SECOND        (drv_usectohz(1000000))

 122 /*
 123  * The data should be "quiet" for this interval (in seconds) before the
 124  * vhci cached data is flushed to the disk.
 125  */
 126 static int mdi_vhcache_flush_delay = 10;
```

```
128 /* number of seconds the vhcache flush daemon will sleep idle before exiting */
129 static int mdi_vhcache_flush_daemon_idle_time = 60;

131 /*
132  * MDI falls back to discovery of all paths when a bus_config_one fails.
133  * The following parameters can be used to tune this operation.
134  *
135  * mdi_path_discovery_boot
136  *      Number of times path discovery will be attempted during early boot.
137  *      Probably there is no reason to ever set this value to greater than one.
138  *
139  * mdi_path_discovery_postboot
140  *      Number of times path discovery will be attempted after early boot.
141  *      Set it to a minimum of two to allow for discovery of iscsi paths which
142  *      may happen very late during booting.
143  *
144  * mdi_path_discovery_interval
145  *      Minimum number of seconds MDI will wait between successive discovery
146  *      of all paths. Set it to -1 to disable discovery of all paths.
147  */
148 static int mdi_path_discovery_boot = 1;
149 static int mdi_path_discovery_postboot = 2;
150 static int mdi_path_discovery_interval = 10;

152 /*
153  * number of seconds the asynchronous configuration thread will sleep idle
154  * before exiting.
155  */
156 static int mdi_async_config_idle_time = 600;

158 static int mdi_bus_config_cache_hash_size = 256;

160 /* turns off multithreaded configuration for certain operations */
161 static int mdi_mtc_off = 0;

163 /*
164  * The "path" to a pathinfo node is identical to the /devices path to a
165  * devinfo node had the device been enumerated under a pHCI instead of
166  * a vHCI.  This pathinfo "path" is associated with a 'path_instance'.
167  * This association persists across create/delete of the pathinfo nodes,
168  * but not across reboot.
169  */
170 static uint_t           mdi_pathmap_instance = 1;        /* 0 -> any path */
171 static int              mdi_pathmap_hash_size = 256;
172 static kmutex_t         mdi_pathmap_mutex;
173 static mod_hash_t       *mdi_pathmap_bypath;             /* "path"->instance */
174 static mod_hash_t       *mdi_pathmap_byinstance;         /* instance->"path" */
175 static mod_hash_t       *mdi_pathmap_sbyinstance;        /* inst->shortpath */

177 /*
178  * MDI component property name/value string definitions
179  */
180 const char              *mdi_component_prop = "mpxio-component";
181 const char              *mdi_component_prop_vhci = "vhci";
182 const char              *mdi_component_prop_phci = "phci";
183 const char              *mdi_component_prop_client = "client";

185 /*
186  * MDI client global unique identifier property name
187  */
188 const char              *mdi_client_guid_prop = "client-guid";

190 /*
191  * MDI client load balancing property name/value string definitions
192  */
193 const char              *mdi_load_balance = "load-balance";
```

```
194 const char              *mdi_load_balance_none = "none";
195 const char              *mdi_load_balance_rr = "round-robin";
196 const char              *mdi_load_balance_lba = "logical-block";

198 /*
199  * Obsolete vHCI class definition; to be removed after Leadville update
200  */
201 const char *mdi_vhci_class_scsi = MDI_HCI_CLASS_SCSI;

203 static char vhci_greeting[] =
204         "\tThere already exists one vHCI driver for class %s\n"
205         "\tOnly one vHCI driver for each class is allowed\n";

207 /*
208  * Static function prototypes
209  */
210 static int              i_mdi_phci_offline(dev_info_t *, uint_t);
211 static int              i_mdi_client_offline(dev_info_t *, uint_t);
212 static int              i_mdi_phci_pre_detach(dev_info_t *, ddi_detach_cmd_t);
213 static void             i_mdi_phci_post_detach(dev_info_t *,
214                         ddi_detach_cmd_t, int);
215 static int              i_mdi_client_pre_detach(dev_info_t *,
216                         ddi_detach_cmd_t);
217 static void             i_mdi_client_post_detach(dev_info_t *,
218                         ddi_detach_cmd_t, int);
219 static void             i_mdi_pm_hold_pip(mdi_pathinfo_t *);
220 static void             i_mdi_pm_rele_pip(mdi_pathinfo_t *);
221 static int              i_mdi_lba_lb(mdi_client_t *ct,
222                         mdi_pathinfo_t **ret_pip, struct buf *buf);
223 static void             i_mdi_pm_hold_client(mdi_client_t *, int);
224 static void             i_mdi_pm_rele_client(mdi_client_t *, int);
225 static void             i_mdi_pm_reset_client(mdi_client_t *);
226 static int              i_mdi_power_all_phci(mdi_client_t *);
227 static void             i_mdi_log_sysevent(dev_info_t *, char *, char *);

230 /*
231  * Internal mdi_pathinfo node functions
232  */
233 static void             i_mdi_pi_kstat_destroy(mdi_pathinfo_t *);

235 static mdi_vhci_t       *i_mdi_vhci_class2vhci(char *);
236 static mdi_vhci_t       *i_devi_get_vhci(dev_info_t *);
237 static mdi_phci_t       *i_devi_get_phci(dev_info_t *);
238 static void             i_mdi_phci_lock(mdi_phci_t *, mdi_pathinfo_t *);
239 static void             i_mdi_phci_unlock(mdi_phci_t *);
240 static mdi_pathinfo_t   *i_mdi_pi_alloc(mdi_phci_t *, char *, mdi_client_t *);
241 static void             i_mdi_phci_add_path(mdi_phci_t *, mdi_pathinfo_t *);
242 static void             i_mdi_client_add_path(mdi_client_t *, mdi_pathinfo_t *);
243 static void             i_mdi_pi_free(mdi_phci_t *ph, mdi_pathinfo_t *,
244                         mdi_client_t *);
245 static void             i_mdi_phci_remove_path(mdi_phci_t *, mdi_pathinfo_t *);
246 static void             i_mdi_client_remove_path(mdi_client_t *,
247                         mdi_pathinfo_t *);

249 static int              i_mdi_pi_state_change(mdi_pathinfo_t *,
250                         mdi_pathinfo_state_t, int);
251 static int              i_mdi_pi_offline(mdi_pathinfo_t *, int);
252 static dev_info_t       *i_mdi_devinfo_create(mdi_vhci_t *, char *, char *,
253                         char **, int);
254 static dev_info_t       *i_mdi_devinfo_find(mdi_vhci_t *, char *, char *);
255 static int              i_mdi_devinfo_remove(dev_info_t *, dev_info_t *, int);
256 static int              i_mdi_is_child_present(dev_info_t *, dev_info_t *);
257 static mdi_client_t     *i_mdi_client_alloc(mdi_vhci_t *, char *, char *);
258 static void             i_mdi_client_enlist_table(mdi_vhci_t *, mdi_client_t *);
259 static void             i_mdi_client_delist_table(mdi_vhci_t *, mdi_client_t *);
```

```
260 static mdi_client_t        *i_mdi_client_find(mdi_vhci_t *, char *, char *);
261 static void                i_mdi_client_update_state(mdi_client_t *);
262 static int                 i_mdi_client_compute_state(mdi_client_t *,
263                                    mdi_phci_t *);
264 static void                i_mdi_client_lock(mdi_client_t *, mdi_pathinfo_t *);
265 static void                i_mdi_client_unlock(mdi_client_t *);
266 static int                 i_mdi_client_free(mdi_vhci_t *, mdi_client_t *);
267 static mdi_client_t        *i_devi_get_client(dev_info_t *);
268 /*
269  * NOTE: this will be removed once the NWS files are changed to use the new
270  * mdi_{enable,disable}_path interfaces
271  */
272 static int                 i_mdi_pi_enable_disable(dev_info_t *, dev_info_t *,
273                                    int, int);
274 static mdi_pathinfo_t      *i_mdi_enable_disable_path(mdi_pathinfo_t *pip,
275                                    mdi_vhci_t *vh, int flags, int op);
276 /*
277  * Failover related function prototypes
278  */
279 static int                 i_mdi_failover(void *);

281 /*
282  * misc internal functions
283  */
284 static int                 i_mdi_get_hash_key(char *);
285 static int                 i_map_nvlist_error_to_mdi(int);
286 static void                i_mdi_report_path_state(mdi_client_t *,
287                                    mdi_pathinfo_t *);

289 static void                setup_vhci_cache(mdi_vhci_t *);
290 static int                 destroy_vhci_cache(mdi_vhci_t *);
291 static int                 stop_vhcache_async_threads(mdi_vhci_config_t *);
292 static boolean_t           stop_vhcache_flush_thread(void *, int);
293 static void                free_string_array(char **, int);
294 static void                free_vhcache_phci(mdi_vhcache_phci_t *);
295 static void                free_vhcache_pathinfo(mdi_vhcache_pathinfo_t *);
296 static void                free_vhcache_client(mdi_vhcache_client_t *);
297 static int                 mainnvl_to_vhcache(mdi_vhci_cache_t *, nvlist_t *);
298 static nvlist_t            *vhcache_to_mainnvl(mdi_vhci_cache_t *);
299 static void                vhcache_phci_add(mdi_vhci_config_t *, mdi_phci_t *);
300 static void                vhcache_phci_remove(mdi_vhci_config_t *, mdi_phci_t *);
301 static void                vhcache_pi_add(mdi_vhci_config_t *,
302                                    struct mdi_pathinfo *);
303 static void                vhcache_pi_remove(mdi_vhci_config_t *,
304                                    struct mdi_pathinfo *);
305 static void                free_phclient_path_list(mdi_phys_path_t *);
306 static void                sort_vhcache_paths(mdi_vhcache_client_t *);
307 static int                 flush_vhcache(mdi_vhci_config_t *, int);
308 static void                vhcache_dirty(mdi_vhci_config_t *);
309 static void                free_async_client_config(mdi_async_client_config_t *);
310 static void                single_threaded_vhconfig_enter(mdi_vhci_config_t *);
311 static void                single_threaded_vhconfig_exit(mdi_vhci_config_t *);
312 static nvlist_t            *read_on_disk_vhci_cache(char *);
313 extern int                 fread_nvlist(char *, nvlist_t **);
314 extern int                 fwrite_nvlist(char *, nvlist_t *);

316 /* called once when first vhci registers with mdi */
317 static void
318 i_mdi_init()
319 {
320        static int initialized = 0;

322        if (initialized)
323               return;
324        initialized = 1;
```

```
326        mutex_init(&mdi_mutex, NULL, MUTEX_DEFAULT, NULL);

328        /* Create our taskq resources */
329        mdi_taskq = taskq_create("mdi_taskq", mdi_taskq_n_threads,
330             MDI_TASKQ_PRI, MDI_TASKQ_MINALLOC, MDI_TASKQ_MAXALLOC,
331             TASKQ_PREPOPULATE | TASKQ_CPR_SAFE);
332        ASSERT(mdi_taskq != NULL);       /* taskq_create never fails */

334        /* Allocate ['path_instance' <-> "path"] maps */
335        mutex_init(&mdi_pathmap_mutex, NULL, MUTEX_DRIVER, NULL);
336        mdi_pathmap_bypath = mod_hash_create_strhash(
337             "mdi_pathmap_bypath", mdi_pathmap_hash_size,
338             mod_hash_null_valdtor);
339        mdi_pathmap_byinstance = mod_hash_create_idhash(
340             "mdi_pathmap_byinstance", mdi_pathmap_hash_size,
341             mod_hash_null_valdtor);
342        mdi_pathmap_sbyinstance = mod_hash_create_idhash(
343             "mdi_pathmap_sbyinstance", mdi_pathmap_hash_size,
344             mod_hash_null_valdtor);
345 }
```
_____**unchanged_portion_omitted_**

```
1439 /*
1440  * i_mdi_client_free():
1441  *            Free a client component
1442  */
1443 static int
1444 i_mdi_client_free(mdi_vhci_t *vh, mdi_client_t *ct)
1445 {
1446        int            rv = MDI_SUCCESS;
1447        int            flags = ct->ct_flags;
1448        dev_info_t     *cdip;
1449        dev_info_t     *vdip;

1451        ASSERT(MDI_VHCI_CLIENT_LOCKED(vh));

1453        vdip = vh->vh_dip;
1454        cdip = ct->ct_dip;

1456        (void) ndi_prop_remove(DDI_DEV_T_NONE, cdip, MDI_CLIENT_GUID_PROP);
1457        DEVI(cdip)->devi_mdi_component &= ~MDI_COMPONENT_CLIENT;
1458        DEVI(cdip)->devi_mdi_client = NULL;

1460        /*
1461         * Clear out back ref. to dev_info_t node
1462         */
1463        ct->ct_dip = NULL;

1465        /*
1466         * Remove this client from our hash queue
1467         */
1468        i_mdi_client_delist_table(vh, ct);

1470        /*
1471         * Uninitialize and free the component
1472         */
1473        kmem_free(ct->ct_drvname, strlen(ct->ct_drvname) + 1);
1474        kmem_free(ct->ct_guid, strlen(ct->ct_guid) + 1);
1475        kmem_free(ct->ct_lb_args, sizeof (client_lb_args_t));
1476        cv_destroy(&ct->ct_failover_cv);
1477        cv_destroy(&ct->ct_unstable_cv);
1478        cv_destroy(&ct->ct_powerchange_cv);
1479        mutex_destroy(&ct->ct_mutex);
1480        kmem_free(ct, sizeof (*ct));
```

```
1481            if (cdip != NULL) {
1482                MDI_VHCI_CLIENT_UNLOCK(vh);
1483                (void) i_mdi_devinfo_remove(vdip, cdip, flags);
1484                MDI_VHCI_CLIENT_LOCK(vh);

1485            }
1486            return (rv);
1487 }
_____unchanged_portion_omitted_
```