

new/usr/src/uts/common/io/comstar/port/iscsit/iscsit.c

1

```
*****
94820 Mon May 5 11:11:12 2014
new/usr/src/uts/common/io/comstar/port/iscsit/iscsit.c
4780 comstar iSCSI target shouldn't abuse ddi_get_time(9f)
Reviewed by: Eric Diven <eric.diven@delphix.com>
Reviewed by: Robert Mustacchi <rm@joyent.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) 2008, 2010, Oracle and/or its affiliates. All rights reserved.
23 *
24 * Copyright 2014 Nexenta Systems, Inc. All rights reserved.
24 * Copyright 2013 Nexenta Systems, Inc. All rights reserved.
25 */

27 #include <sys/cpuvar.h>
28 #include <sys/types.h>
29 #include <sys/conf.h>
30 #include <sys/stat.h>
31 #include <sys/file.h>
32 #include <sys/ddi.h>
33 #include <sys/sunddi.h>
34 #include <sys/modctl.h>
35 #include <sys/sysmacros.h>
36 #include <sys/socket.h>
37 #include <sys/strsubr.h>
38 #include <sys/nvpair.h>

40 #include <sys/stmf.h>
41 #include <sys/stmf_ioctl.h>
42 #include <sys/portif.h>
43 #include <sys/idm/idm.h>
44 #include <sys/idm/idm_conn_sm.h>

46 #include "iscsit_isns.h"
47 #include "iscsit.h"

49 #define ISCSIT_VERSION BUILD_DATE "-1.18dev"
50 #define ISCSIT_NAME_VERSION "COMSTAR ISCSIT v" ISCSIT_VERSION

52 /*
53  * DDI entry points.
54  */
55 static int iscsit_drv_attach(dev_info_t *, ddi_attach_cmd_t);
56 static int iscsit_drv_detach(dev_info_t *, ddi_detach_cmd_t);
57 static int iscsit_drv_getinfo(dev_info_t *, ddi_info_cmd_t, void *, void **);
58 static int iscsit_drv_open(dev_t *, int, int, cred_t *);
```

new/usr/src/uts/common/io/comstar/port/iscsit/iscsit.c

2

```
59 static int iscsit_drv_close(dev_t, int, int, cred_t *);
60 static boolean_t iscsit_drv_busy(void);
61 static int iscsit_drv_ioctl(dev_t, int, intp_t, int, cred_t *, int *);

63 extern struct mod_ops mod_miscops;

66 static struct cb_ops iscsit_cb_ops = {
67     iscsit_drv_open,          /* cb_open */
68     iscsit_drv_close,        /* cb_close */
69     nodev,                    /* cb_strategy */
70     nodev,                    /* cb_print */
71     nodev,                    /* cb_dump */
72     nodev,                    /* cb_read */
73     nodev,                    /* cb_write */
74     iscsit_drv_ioctl,        /* cb_ioctl */
75     nodev,                    /* cb_devmap */
76     nodev,                    /* cb_mmap */
77     nodev,                    /* cb_segmap */
78     nochpoll,                /* cb_chpoll */
79     ddi_prop_op,             /* cb_prop_op */
80     NULL,                    /* cb_streamtab */
81     D_MP,                    /* cb_flag */
82     CB_REV,                  /* cb_rev */
83     nodev,                    /* cb_aread */
84     nodev,                    /* cb_awrite */
85 };
    unchanged_portion_omitted

3093 /*
3094  * iscsit_check_cmdsnd_and_queue
3095  *
3096  * Independent of the order in which the iSCSI target receives non-immediate
3097  * command PDU across the entire session and any multiple connections within
3098  * the session, the target must deliver the commands to the SCSI layer in
3099  * CmdSN order. So out-of-order non-immediate commands are queued up on a
3100  * session-wide wait queue. Duplicate commands are ignored.
3101  *
3102  */
3103 static int
3104 iscsit_check_cmdsnd_and_queue(idm_pdu_t *rx_pdu)
3105 {
3106     idm_conn_t          *ic = rx_pdu->isp_ic;
3107     iscsit_conn_t       *ict = ic->ic_handle;
3108     iscsit_sess_t       *ist = ict->ict_sess;
3109     iscsi_scsi_cmd_hdr_t *hdr = (iscsi_scsi_cmd_hdr_t *)rx_pdu->isp_hdr;

3111     mutex_enter(&ist->ist_sn_mutex);
3112     if (hdr->opcode & ISCSI_OP_IMMEDIATE) {
3113         /* do not queue, handle it immediately */
3114         DTRACE_PROBE2(immediate_cmd, iscsit_sess_t *, ist,
3115             idm_pdu_t *, rx_pdu);
3116         mutex_exit(&ist->ist_sn_mutex);
3117         return (ISCSIT_CMDSND_EQ_EXPCMDSND);
3118     }
3119     if (iscsit_sna_lt(ist->ist_expcmdsn, ntohl(hdr->cmdsnd))) {
3120         /*
3121          * Out-of-order commands (cmdSN higher than ExpCmdSN)
3122          * are staged on a fixed-size circular buffer until
3123          * the missing command is delivered to the SCSI layer.
3124          * Irrespective of the order of insertion into the
3125          * staging queue, the commands are processed out of the
3126          * queue in cmdSN order only.
3127          */
3128         rx_pdu->isp_queue_time = gethrtime();
3128         rx_pdu->isp_queue_time = ddi_get_time();
```

```

3129         iscsit_add_pdu_to_queue(ist, rx_pdu);
3130         mutex_exit(&ist->ist_sn_mutex);
3131         return (ISCSIT_CMDSN_GT_EXPCMDSN);
3132     } else if (iscsit_sna_lt(ntohl(hdr->cmds), ist->ist_expcmdsn)) {
3133         DTRACE_PROBE3(cmdsn_lt_expcmdsn, iscsit_sess_t *, ist,
3134             iscsit_conn_t *, ict, idm_pdu_t *, rx_pdu);
3135         mutex_exit(&ist->ist_sn_mutex);
3136         return (ISCSIT_CMDSN_LT_EXPCMDSN);
3137     } else {
3138         mutex_exit(&ist->ist_sn_mutex);
3139         return (ISCSIT_CMDSN_EQ_EXPCMDSN);
3140     }
3141 }

```

unchanged portion omitted

```

3356 static void
3357 iscsit_rxpdu_queue_monitor_session(iscsit_sess_t *ist)
3358 {
3359     iscsit_cbuf_t *cbuf = ist->ist_rxpdu_queue;
3360     idm_pdu_t *next_pdu = NULL;
3361     uint32_t index, next_cmdsn, i;
3362
3363     /*
3364     * Assume that all PDUs in the staging queue have a cmds >= expcmds.
3365     * Starting with the expcmds, iterate over the staged PDUs to find
3366     * the next PDU with a wait time greater than the threshold. If found
3367     * advance the staged PDU to the SCSI layer, skipping over the missing
3368     * PDU(s) to get past the hole in the command sequence. It is up to
3369     * the initiator to note that the target has not acknowledged a cmds
3370     * and take appropriate action.
3371     *
3372     * Since the PDU(s) arrive in any random order, it is possible that
3373     * that the actual wait time for a particular PDU is much longer than
3374     * the defined threshold. e.g. Consider a case where commands are sent
3375     * over 4 different connections, and cmds = 1004 arrives first, then
3376     * 1003, and 1002 and 1001 are lost due to a connection failure.
3377     * So now 1003 is waiting for 1002 to be delivered, and although the
3378     * wait time of 1004 > wait time of 1003, only 1003 will be considered
3379     * by the monitor thread. 1004 will be automatically processed by
3380     * iscsit_process_pdu_in_queue() once the scan is complete and the
3381     * expcmds becomes current.
3382     */
3383     mutex_enter(&ist->ist_sn_mutex);
3384     cbuf = ist->ist_rxpdu_queue;
3385     if (cbuf->cb_num_elems == 0) {
3386         mutex_exit(&ist->ist_sn_mutex);
3387         return;
3388     }
3389     for (next_pdu = NULL, i = 0; i < cbuf->cb_num_elems; i++) {
3390         next_cmdsn = ist->ist_expcmdsn + i; /* start at expcmds */
3391         index = next_cmdsn % ISCSIT_RXPDU_QUEUE_LEN;
3392         if ((next_pdu = cbuf->cb_buffer[index]) != NULL) {
3393             /*
3394             * If the PDU wait time has not exceeded threshold
3395             * stop scanning the staging queue until the timer
3396             * fires again
3397             */
3398             if ((gethrtime() - next_pdu->isp_queue_time)
3399                 < (rxpdu_queue_threshold * NANOS)) {
3400                 if ((ddi_get_time() - next_pdu->isp_queue_time)
3401                     < rxpdu_queue_threshold) {
3402                     mutex_exit(&ist->ist_sn_mutex);
3403                     return;
3404                 }
3405             }
3406             /*
3407             * Remove the next PDU from the queue and post it

```

```

3405         * to the SCSI layer, skipping over the missing
3406         * PDU. Stop scanning the staging queue until
3407         * the monitor timer fires again
3408         */
3409         (void) iscsit_remove_pdu_from_queue(ist, next_cmdsn);
3410         mutex_exit(&ist->ist_sn_mutex);
3411         DTRACE_PROBE3(advanced_to_blocked_cmds,
3412             iscsit_sess_t *, ist, idm_pdu_t *, next_pdu,
3413             uint32_t, next_cmdsn);
3414         iscsit_post_staged_pdu(next_pdu);
3415         /* Deliver any subsequent PDUs immediately */
3416         iscsit_process_pdu_in_queue(ist);
3417         return;
3418     }
3419     /*
3420     * Skipping over i PDUs, e.g. a case where commands 1001 and
3421     * 1002 are lost in the network, skip over both and post 1003
3422     * expcmds then becomes 1004 at the end of the scan.
3423     */
3424     DTRACE_PROBE2(skipping_over_cmds, iscsit_sess_t *, ist,
3425         uint32_t, next_cmdsn);
3426 }
3427 /*
3428 * following the assumption, staged cmds >= expcmds, this statement
3429 * is never reached.
3430 */
3431 }

```

unchanged portion omitted

```

*****
14701 Mon May 5 11:11:12 2014
new/usr/src/uts/common/sys/idm/idm_impl.h
4780 comstar iSCSI target shouldn't abuse ddi_get_time(9f)
Reviewed by: Eric Diven <eric.diven@delphix.com>
Reviewed by: Robert Mustacchi <rm@joyent.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) 2008, 2010, Oracle and/or its affiliates. All rights reserved.
23 */
24 /*
25  * Copyright 2014 Nexenta Systems, Inc. All rights reserved.
26  * Copyright 2011 Nexenta Systems, Inc. All rights reserved.
27 */

28 #ifndef _IDM_IMPL_H_
29 #define _IDM_IMPL_H_

30 #ifdef __cplusplus
31 extern "C" {
32 #endif

33 #include <sys/avl.h>
34 #include <sys/socket_impl.h>
35 #include <sys/taskq_impl.h>

36 /*
37  * IDM lock order:
38  *
39  * idm_taskid_table_lock, idm_task_t.idt_mutex
40 */

41 #define CF_LOGIN_READY      0x00000001
42 #define CF_INITIAL_LOGIN   0x00000002
43 #define CF_ERROR           0x80000000

44 typedef enum {
45     CONN_TYPE_INI = 1,
46     CONN_TYPE_TGT
47 } idm_conn_type_t;
48
49 unchanged_portion_omitted

334 #define PDU_MAX_IOVLEN 12
335 #define IDM_PDU_MAGIC 0x49504455 /* "IPDU" */

336 typedef struct idm_pdu_s {

```

```

338     uint32_t      isp_magic; /* "IPDU" */

339     /*
340      * Internal - Order is vital. idm_tx_link *must* be the second
341      * element in this structure for proper TX PDU ordering.
342      */
343     list_node_t   idm_tx_link;

344     list_node_t   isp_client_lnd;

345     idm_conn_t    *isp_ic; /* Must be set */
346     iscsi_hdr_t   *isp_hdr;
347     uint_t        isp_hdrlen;
348     uint8_t       *isp_data;
349     uint_t        isp_dataalen;

350     /* Transport header */
351     void          *isp_transport_hdr;
352     uint32_t      isp_transport_hdrlen;
353     void          *isp_transport_private;

354     /*
355      * isp_data is used for sending SCSI status, NOP, text, scsi and
356      * non-scsi data. Data is received using isp_iov and isp_iovlen
357      * to support data over multiple buffers.
358      */
359     void          *isp_private;
360     idm_pdu_cb_t  *isp_callback;
361     idm_status_t  isp_status;

362     /*
363      * The following four elements are only used in
364      * idm_sorecv_scsideata() currently.
365      */
366     struct iovec  isp_iov[PDU_MAX_IOVLEN];
367     int           isp_iovlen;
368     idm_buf_t     isp_sorx_buf;

369     /* Implementation data for idm_pdu_alloc and sorx PDU cache */
370     uint32_t      isp_flags;
371     uint_t        isp_hdrbuflen;
372     uint_t        isp_databuflen;
373     hrtime_t      isp_queue_time;
374     time_t        isp_queue_time;

375     /* Taskq dispatching state for deferred PDU */
376     taskq_ent_t  isp_tqent;
377 } idm_pdu_t;
378
379 unchanged_portion_omitted

```