

new/usr/src/pkg/manifests/system-test-ostest.mf

1

```
*****
3734 Thu Jun 14 17:15:57 2018
new/usr/src/pkg/manifests/system-test-ostest.mf
9600 LDT still not happy under KPTI
*****
1 #
2 # This file and its contents are supplied under the terms of the
3 # Common Development and Distribution License ("CDDL"), version 1.0.
4 # You may only use this file in accordance with the terms of version
5 # 1.0 of the CDDL.
6 #
7 # A full copy of the text of the CDDL should have accompanied this
8 # source. A copy of the CDDL is also available via the Internet at
9 # http://www.illumos.org/license/CDDL.
10 #
11 #
12 #
13 # Copyright (c) 2012, 2016 by Delphix. All rights reserved.
14 # Copyright 2014, OmniTI Computer Consulting, Inc. All rights reserved.
15 # Copyright 2018 Joyent, Inc.
16 # Copyright 2017 Joyent, Inc.
17 #
18 set name=pkg.fmri value=pkg:/system/test/ostest@$(PKGVERS)
19 set name=pkg.description value="Miscellaneous OS Unit Tests"
20 set name=pkg.summary value="OS Unit Test Suite"
21 set name=info.classification \
22     value=org.opensolaris.category.2008:Development/System
23 set name=variant.arch value=$(ARCH)
24 dir path=opt/os-tests
25 dir path=opt/os-tests/bin
26 dir path=opt/os-tests/runfiles
27 dir path=opt/os-tests/tests
28 dir path=opt/os-tests/tests/file-locking
29 $(i386_ONLY)dir path=opt/os-tests/tests/i386
30 dir path=opt/os-tests/tests/pf_key
31 dir path=opt/os-tests/tests/sdevfs
32 dir path=opt/os-tests/tests/secflags
33 dir path=opt/os-tests/tests/sigqueue
34 dir path=opt/os-tests/tests/sockfs
35 dir path=opt/os-tests/tests/stress
36 file path=opt/os-tests/README mode=0444
37 file path=opt/os-tests/bin/ostest mode=0555
38 file path=opt/os-tests/runfiles/default.run mode=0444
39 file path=opt/os-tests/tests/epoll_test mode=0555
40 file path=opt/os-tests/tests/file-locking/acquire-lock.32 mode=0555
41 file path=opt/os-tests/tests/file-locking/acquire-lock.64 mode=0555
42 file path=opt/os-tests/tests/file-locking/runtests.32 mode=0555
43 file path=opt/os-tests/tests/file-locking/runtests.64 mode=0555
44 $(i386_ONLY)file path=opt/os-tests/tests/i386/ldt mode=0555
45 file path=opt/os-tests/tests/pf_key/acquire-compare mode=0555
46 file path=opt/os-tests/tests/pf_key/acquire-spray mode=0555
47 file path=opt/os-tests/tests/pf_key/eacq-enabler mode=0555
48 file path=opt/os-tests/tests/pf_key/kmc-update mode=0555
49 file path=opt/os-tests/tests/pf_key/kmc-updater mode=0555
50 file path=opt/os-tests/tests/poll_test mode=0555
51 file path=opt/os-tests/tests/sdevfs/sdevfs_eisdir mode=0555
52 file path=opt/os-tests/tests/secflags/addr32 mode=0555
53 file path=opt/os-tests/tests/secflags/addr64 mode=0555
54 file path=opt/os-tests/tests/secflags/secflags_aslr mode=0555
55 file path=opt/os-tests/tests/secflags/secflags_core mode=0555
56 file path=opt/os-tests/tests/secflags/secflags_dts mode=0555
57 file path=opt/os-tests/tests/secflags/secflags_elfdump mode=0555
58 file path=opt/os-tests/tests/secflags/secflags_forbidnullmap mode=0555
59 file path=opt/os-tests/tests/secflags/secflags_limits mode=0555
60 file path=opt/os-tests/tests/secflags/secflags_noexecstack mode=0555
```

new/usr/src/pkg/manifests/system-test-ostest.mf

2

```
61 file path=opt/os-tests/tests/secflags/secflags_proc mode=0555
62 file path=opt/os-tests/tests/secflags/secflags_psecflags mode=0555
63 file path=opt/os-tests/tests/secflags/secflags_syscall mode=0555
64 file path=opt/os-tests/tests/secflags/secflags_truss mode=0555
65 file path=opt/os-tests/tests/secflags/secflags_zonecfg mode=0555
66 file path=opt/os-tests/tests/secflags/stacky mode=0555
67 file path=opt/os-tests/tests/sigqueue/sigqueue_queue_size mode=0555
68 file path=opt/os-tests/tests/sockfs/conn mode=0555
69 file path=opt/os-tests/tests/sockfs/dgram mode=0555
70 file path=opt/os-tests/tests/sockfs/drop_priv mode=0555
71 file path=opt/os-tests/tests/sockfs/nosignal mode=0555
72 file path=opt/os-tests/tests/sockfs/sockpair mode=0555
73 file path=opt/os-tests/tests/spoof-ras mode=0555
74 file path=opt/os-tests/tests/stress/dladm-kstat mode=0555
75 license cr_Sun license=cr_Sun
76 license lic_CDDL license=lic_CDDL
77 depend fmri=pkg:/network/telnet type=require
78 depend fmri=system/test/testrunner type=require
```

new/usr/src/test/os-tests/runfiles/default.run

1

```
*****
1595 Thu Jun 14 17:15:57 2018
new/usr/src/test/os-tests/runfiles/default.run
9600 LDT still not happy under KPTI
*****
1 #
2 # This file and its contents are supplied under the terms of the
3 # Common Development and Distribution License ("CDDL"), version 1.0.
4 # You may only use this file in accordance with the terms of version
5 # 1.0 of the CDDL.
6 #
7 # A full copy of the text of the CDDL should have accompanied this
8 # source. A copy of the CDDL is also available via the Internet at
9 # http://www.illumos.org/license/CDDL.
10 #
11 #
12 #
13 # Copyright (c) 2012 by Delphix. All rights reserved.
14 # Copyright 2018 Joyent, Inc.
15 # Copyright 2017 Joyent, Inc.
16 #
17 [DEFAULT]
18 pre =
19 verbose = False
20 quiet = False
21 timeout = 60
22 post =
23 outputdir = /var/tmp/test_results
24 #
25 [/opt/os-tests/tests/poll_test]
26 user = root
27 tests = ['poll_test', 'epoll_test']
28 #
29 [/opt/os-tests/tests/secflags]
30 user = root
31 tests = ['secflags_aslr',
32          'secflags_core',
33          'secflags_dts',
34          'secflags_elfdump',
35          'secflags_forbidnullmap',
36          'secflags_limits',
37          'secflags_noexecstack',
38          'secflags_proc',
39          'secflags_psecflags',
40          'secflags_syscall',
41          'secflags_truss',
42          'secflags_zonecfg']
43 #
44 [/opt/os-tests/tests/sigqueue]
45 tests = ['sigqueue_queue_size']
46 #
47 [/opt/os-tests/tests/sdevfs]
48 user = root
49 tests = ['sdevfs_eisdir']
50 #
51 [/opt/os-tests/tests/stress]
52 user = root
53 tests = ['dladm-kstat']
54 #
55 [/opt/os-tests/tests/file-locking]
56 tests = ['runtests.32', 'runtests.64']
57 #
58 [/opt/os-tests/tests/sockfs]
59 user = root
60 tests = ['conn', 'dgram', 'drop_priv', 'nosignal', 'sockpair']
```

new/usr/src/test/os-tests/runfiles/default.run

2

```
62 [/opt/os-tests/tests/pf_key]
63 user = root
64 tests = ['acquire-compare', 'acquire-spray']
65 #
66 [/opt/os-tests/tests/i386]
67 user = root
68 arch = i86pc
69 tests = ['ldt']
```

new/usr/src/test/os-tests/tests/Makefile

1

```
*****
659 Thu Jun 14 17:15:58 2018
new/usr/src/test/os-tests/tests/Makefile
9600 LDT still not happy under KPTI
*****
1 #
2 # This file and its contents are supplied under the terms of the
3 # Common Development and Distribution License ("CDDL"), version 1.0.
4 # You may only use this file in accordance with the terms of version
5 # 1.0 of the CDDL.
6 #
7 # A full copy of the text of the CDDL should have accompanied this
8 # source. A copy of the CDDL is also available via the Internet at
9 # http://www.illumos.org/license/CDDL.
10 #
11 #
12 #
13 # Copyright (c) 2012, 2016 by Delphix. All rights reserved.
14 # Copyright 2018 Joyent, Inc.
14 # Copyright 2017 Joyent, Inc.
15 #
16 #
17 SUBDIRS_i386 = i386
18
19 SUBDIRS = poll secflags sigqueue spoof-ras sdevfs sockfs stress file-locking \
20         pf_key $(SUBDIRS_$(MACH))
21         pf_key
22
23 include $(SRC)/test/Makefile.com
```

new/usr/src/test/os-tests/tests/i386/Makefile

1

\*\*\*\*\*

833 Thu Jun 14 17:15:58 2018

new/usr/src/test/os-tests/tests/i386/Makefile

9600 LDT still not happy under KPTI

\*\*\*\*\*

```
1 #
2 # This file and its contents are supplied under the terms of the
3 # Common Development and Distribution License ("CDDL"), version 1.0.
4 # You may only use this file in accordance with the terms of version
5 # 1.0 of the CDDL.
6 #
7 # A full copy of the text of the CDDL should have accompanied this
8 # source. A copy of the CDDL is also available via the Internet at
9 # http://www.illumos.org/license/CDDL.
10 #
11 #
12 #
13 # Copyright 2018 Joyent, Inc.
14 #
15 #
16 include $(SRC)/cmd/Makefile.cmd
17 include $(SRC)/test/Makefile.com
18 #
19 PROG += ldt
20 #
21 ROOTOPTPKG = $(ROOT)/opt/os-tests
22 TESTDIR = $(ROOTOPTPKG)/tests/i386
23 #
24 CSTD = $(CSTD_GNU99)
25 #
26 CMDS = $(PROG:%=$(TESTDIR)/%)
27 $(CMDS) := FILEMODE = 0555
28 #
29 all: $(PROG)
30 #
31 install: all $(CMDS)
32 #
33 lint:
34 #
35 clobber: clean
36     -$(RM) $(PROG)
37 #
38 clean:
39 #
40 $(CMDS): $(TESTDIR) $(PROG)
41 #
42 $(TESTDIR):
43     $(INS.dir)
44 #
45 $(TESTDIR)/%: %
46     $(INS.file)
```

new/usr/src/test/os-tests/tests/i386/ldt.c

1

```
*****
1672 Thu Jun 14 17:15:58 2018
new/usr/src/test/os-tests/tests/i386/ldt.c
9600 LDT still not happy under KPTI
*****
1 /*
2  * This file and its contents are supplied under the terms of the
3  * Common Development and Distribution License ("CDDL"), version 1.0.
4  * You may only use this file in accordance with the terms of version
5  * 1.0 of the CDDL.
6  *
7  * A full copy of the text of the CDDL should have accompanied this
8  * source. A copy of the CDDL is also available via the Internet at
9  * http://www.illumos.org/license/CDDL.
10 */

12 /*
13  * Copyright 2018 Joyent, Inc.
14 */

16 #include <sys/types.h>
17 #include <sys/sysi86.h>
18 #include <sys/segments.h>
19 #include <sys/segment.h>
20 #include <unistd.h>
21 #include <string.h>
22 #include <errno.h>
23 #include <pthread.h>
24 #include <err.h>

26 char foo[4096];

28 static void *
29 donothing(void *nothing)
30 {
31     sleep(5);
32     return (NULL);
33 }

35 int
36 main(void)
37 {
38     pthread_t tid;

40     /*
41      * This first is similar to what sbcl does in some variants. Note the
42      * SDT_MEMRW (not SDT_MEMRWA) so we check that the kernel is forcing the
43      * 'accessed' bit too.
44      */
45     int sel = SEL_LDT(7);

47     struct ssd ssd = { sel, (unsigned long)&foo, 4096,
48         SDT_MEMRW | (SEL_UPL << 5) | (1 << 7), 0x4 };

50     if (sysi86(SI86DSCR, &ssd) < 0)
51         err(-1, "failed to setup segment");

53     __asm__ __volatile__ ("mov %0, %%fs" : : "r" (sel));

55     ssd.accl = 0;

57     if (sysi86(SI86DSCR, &ssd) == 0)
58         errx(-1, "removed in-use segment?");

60     __asm__ __volatile__ ("mov %0, %%fs" : : "r" (0));
```

new/usr/src/test/os-tests/tests/i386/ldt.c

2

```
62     if (sysi86(SI86DSCR, &ssd) < 0)
63         err(-1, "failed to remove segment");

65     for (int i = 0; i < MAXNLDT; i++) {
66         ssd.sel = SEL_LDT(i);
67         (void) sysi86(SI86DSCR, &ssd);
68     }

70     for (int i = 0; i < 10; i++)
71         pthread_create(&tid, NULL, donothing, NULL);

73     if (forkall() == 0) {
74         sleep(2);
75         _exit(0);
76     }

78     sleep(6);
79     return (0);
80 }
```

new/usr/src/test/test-runner/cmd/run

1

```
*****
32392 Thu Jun 14 17:15:58 2018
new/usr/src/test/test-runner/cmd/run
9600 LDT still not happy under KPTI
*****
1 #!@PYTHON@

3 #
4 # This file and its contents are supplied under the terms of the
5 # Common Development and Distribution License ("CDDL"), version 1.0.
6 # You may only use this file in accordance with the terms of version
7 # 1.0 of the CDDL.
8 #
9 # A full copy of the text of the CDDL should have accompanied this
10 # source. A copy of the CDDL is also available via the Internet at
11 # http://www.illumos.org/license/CDDL.
12 #

14 #
15 # Copyright (c) 2012, 2016 by Delphix. All rights reserved.
16 # Copyright (c) 2017, Chris Fraire <cfraire@me.com>.
17 # Copyright 2018 Joyent, Inc.
18 #

20 import ConfigParser
21 import os
22 import logging
23 import platform
24 from logging.handlers import WatchedFileHandler
25 from datetime import datetime
26 from optparse import OptionParser
27 from pwd import getpwnam
28 from pwd import getpwuid
29 from select import select
30 from subprocess import PIPE
31 from subprocess import Popen
32 from sys import argv
33 from sys import maxint
34 from threading import Timer
35 from time import time

37 BASEDIR = '/var/tmp/test_results'
38 KILL = '/usr/bin/kill'
39 TRUE = '/usr/bin/true'
40 SUDO = '/usr/bin/sudo'

42 # Custom class to reopen the log file in case it is forcibly closed by a test.
43 class WatchedFileHandlerClosed(WatchedFileHandler):
44     """Watch files, including closed files.
45     Similar to (and inherits from) logging.handler.WatchedFileHandler,
46     except that IOErrors are handled by reopening the stream and retrying.
47     This will be retried up to a configurable number of times before
48     giving up, default 5.
49     """

51     def __init__(self, filename, mode='a', encoding=None, delay=0, max_tries=5):
52         self.max_tries = max_tries
53         self.tries = 0
54         WatchedFileHandler.__init__(self, filename, mode, encoding, delay)

56     def emit(self, record):
57         while True:
58             try:
59                 WatchedFileHandler.emit(self, record)
60                 self.tries = 0
61                 return
```

new/usr/src/test/test-runner/cmd/run

2

```
62         except IOError as err:
63             if self.tries == self.max_tries:
64                 raise
65             self.stream.close()
66             self.stream = self._open()
67             self.tries += 1

69 class Result(object):
70     total = 0
71     runresults = {'PASS': 0, 'FAIL': 0, 'SKIP': 0, 'KILLED': 0}

73     def __init__(self):
74         self.starttime = None
75         self.returncode = None
76         self.runtime = ''
77         self.stdout = []
78         self.stderr = []
79         self.result = ''

81     def done(self, proc, killed):
82         """
83         Finalize the results of this Cmd.
84         """
85         Result.total += 1
86         m, s = divmod(time() - self.starttime, 60)
87         self.runtime = '%02d:%02d' % (m, s)
88         self.returncode = proc.returncode
89         if killed:
90             self.result = 'KILLED'
91             Result.runresults['KILLED'] += 1
92         elif self.returncode is 0:
93             self.result = 'PASS'
94             Result.runresults['PASS'] += 1
95         elif self.returncode is not 0:
96             self.result = 'FAIL'
97             Result.runresults['FAIL'] += 1

100 class Output(object):
101     """
102     This class is a slightly modified version of the 'Stream' class found
103     here: http://goo.gl/aSGfv
104     """
105     def __init__(self, stream):
106         self.stream = stream
107         self._buf = ''
108         self.lines = []

110     def fileno(self):
111         return self.stream.fileno()

113     def read(self, drain=0):
114         """
115         Read from the file descriptor. If 'drain' set, read until EOF.
116         """
117         while self._read() is not None:
118             if not drain:
119                 break

121     def _read(self):
122         """
123         Read up to 4k of data from this output stream. Collect the output
124         up to the last newline, and append it to any leftover data from a
125         previous call. The lines are stored as a (timestamp, data) tuple
126         for easy sorting/merging later.
127         """
```

```

128     fd = self.fileno()
129     buf = os.read(fd, 4096)
130     if not buf:
131         return None
132     if '\n' not in buf:
133         self._buf += buf
134     return []

136     buf = self._buf + buf
137     tmp, rest = buf.rsplit('\n', 1)
138     self._buf = rest
139     now = datetime.now()
140     rows = tmp.split('\n')
141     self.lines += [(now, r) for r in rows]

144 class Cmd(object):
145     verified_users = []

147     def __init__(self, pathname, outputdir=None, timeout=None, user=None):
148         self.pathname = pathname
149         self.outputdir = outputdir or 'BASEDIR'
150         self.timeout = timeout
151         self.user = user or ''
152         self.killed = False
153         self.result = Result()

155         if self.timeout is None:
156             self.timeout = 60

158     def __str__(self):
159         return "Pathname: %s\nOutputdir: %s\nTimeout: %s\nUser: %s\n" % \
160             (self.pathname, self.outputdir, self.timeout, self.user)

162     def kill_cmd(self, proc):
163         """
164         Kill a running command due to timeout, or ^C from the keyboard. If
165         sudo is required, this user was verified previously.
166         """
167         self.killed = True
168         do_sudo = len(self.user) != 0
169         signal = '-TERM'

171         cmd = [SUDO, KILL, signal, str(proc.pid)]
172         if not do_sudo:
173             del cmd[0]

175         try:
176             kp = Popen(cmd)
177             kp.wait()
178         except:
179             pass

181     def update_cmd_privs(self, cmd, user):
182         """
183         If a user has been specified to run this Cmd and we're not already
184         running as that user, prepend the appropriate sudo command to run
185         as that user.
186         """
187         me = getpwuid(os.getuid())

189         if not user or user is me:
190             return cmd

192         ret = '%s -E -u %s %s' % (SUDO, user, cmd)
193         return ret.split(' ')

```

```

195     def collect_output(self, proc):
196         """
197         Read from stdout/stderr as data becomes available, until the
198         process is no longer running. Return the lines from the stdout and
199         stderr Output objects.
200         """
201         out = Output(proc.stdout)
202         err = Output(proc.stderr)
203         res = []
204         while proc.returncode is None:
205             proc.poll()
206             res = select([out, err], [], [], .1)
207             for fd in res[0]:
208                 fd.read()
209         for fd in res[0]:
210             fd.read(drain=1)

212         return out.lines, err.lines

214     def run(self, options):
215         """
216         This is the main function that runs each individual test.
217         Determine whether or not the command requires sudo, and modify it
218         if needed. Run the command, and update the result object.
219         """
220         if options.dryrun is True:
221             print self
222             return

224         privcmd = self.update_cmd_privs(self.pathname, self.user)
225         try:
226             old = os.umask(0)
227             if not os.path.isdir(self.outputdir):
228                 os.makedirs(self.outputdir, mode=0777)
229             os.umask(old)
230         except OSError, e:
231             fail('%s' % e)

233         try:
234             self.result.starttime = time()
235             proc = Popen(privcmd, stdout=PIPE, stderr=PIPE, stdin=PIPE)
236             proc.stdin.close()

238             # Allow a special timeout value of 0 to mean infinity
239             if int(self.timeout) == 0:
240                 self.timeout = maxint
241             t = Timer(int(self.timeout), self.kill_cmd, [proc])
242             t.start()
243             self.result.stdout, self.result.stderr = self.collect_output(proc)
244         except KeyboardInterrupt:
245             self.kill_cmd(proc)
246             fail('\nRun terminated at user request.')
247         finally:
248             t.cancel()

250         self.result.done(proc, self.killed)

252     def skip(self):
253         """
254         Initialize enough of the test result that we can log a skipped
255         command.
256         """
257         Result.total += 1
258         Result.runresults['SKIP'] += 1
259         self.result.stdout = self.result.stderr = []

```

```

260     self.result.starttime = time()
261     m, s = divmod(time() - self.result.starttime, 60)
262     self.result.runtime = '%02d:%02d' % (m, s)
263     self.result.result = 'SKIP'

265     def log(self, logger, options):
266         """
267         This function is responsible for writing all output. This includes
268         the console output, the logfile of all results (with timestamped
269         merged stdout and stderr), and for each test, the unmodified
270         stdout/stderr/merged in it's own file.
271         """
272         if logger is None:
273             return

275         logname = getpwnid(os.getuid()).pw_name
276         user = ' (run as %s)' % (self.user if len(self.user) else logname)
277         msga = 'Test: %s%s' % (self.pathname, user)
278         msgb = '[%s] [%s]' % (self.result.runtime, self.result.result)
279         pad = ' ' * (80 - (len(msga) + len(msgb)))

281         # If -q is specified, only print a line for tests that didn't pass.
282         # This means passing tests need to be logged as DEBUG, or the one
283         # line summary will only be printed in the logfile for failures.
284         if not options.quiet:
285             logger.info('%s%s%s' % (msga, pad, msgb))
286         elif self.result.result is not 'PASS':
287             logger.info('%s%s%s' % (msga, pad, msgb))
288         else:
289             logger.debug('%s%s%s' % (msga, pad, msgb))

291         lines = sorted(self.result.stdout + self.result.stderr,
292                       cmp=lambda x, y: cmp(x[0], y[0]))

294         for dt, line in lines:
295             logger.debug('%s %s' % (dt.strftime("%H:%M:%S.%f ")[:11], line))

297         if len(self.result.stdout):
298             with open(os.path.join(self.outputdir, 'stdout'), 'w') as out:
299                 for _, line in self.result.stdout:
300                     os.write(out.fileno(), '%s\n' % line)
301         if len(self.result.stderr):
302             with open(os.path.join(self.outputdir, 'stderr'), 'w') as err:
303                 for _, line in self.result.stderr:
304                     os.write(err.fileno(), '%s\n' % line)
305         if len(self.result.stdout) and len(self.result.stderr):
306             with open(os.path.join(self.outputdir, 'merged'), 'w') as merged:
307                 for _, line in lines:
308                     os.write(merged.fileno(), '%s\n' % line)

311 class Test(Cmd):
312     props = ['outputdir', 'timeout', 'user', 'pre', 'pre_user', 'post',
313            'post_user']

315     def __init__(self, pathname, outputdir=None, timeout=None, user=None,
316                pre=None, pre_user=None, post=None, post_user=None):
317         super(Test, self).__init__(pathname, outputdir, timeout, user)
318         self.pre = pre or ''
319         self.pre_user = pre_user or ''
320         self.post = post or ''
321         self.post_user = post_user or ''

323     def __str__(self):
324         post_user = pre_user = ''
325         if len(self.pre_user):

```

```

326         pre_user = ' (as %s)' % (self.pre_user)
327         if len(self.post_user):
328             post_user = ' (as %s)' % (self.post_user)
329         return "Pathname: %s\nOutputdir: %s\nTimeout: %d\nPre: %s\nPost: " \
330                "%s\nUser: %s\n" % \
331                (self.pathname, self.outputdir, self.timeout, self.pre,
332                 pre_user, self.post, post_user, self.user)

334     def verify(self, logger):
335         """
336         Check the pre/post scripts, user and Test. Omit the Test from this
337         run if there are any problems.
338         """
339         files = [self.pre, self.pathname, self.post]
340         users = [self.pre_user, self.user, self.post_user]

342         for f in [f for f in files if len(f)]:
343             if not verify_file(f):
344                 logger.info("Warning: Test '%s' not added to this run because"
345                             " it failed verification." % f)
346                 return False

348         for user in [user for user in users if len(user)]:
349             if not verify_user(user, logger):
350                 logger.info("Not adding Test '%s' to this run." %
351                             self.pathname)
352                 return False

354         return True

356     def run(self, logger, options):
357         """
358         Create Cmd instances for the pre/post scripts. If the pre script
359         doesn't pass, skip this Test. Run the post script regardless.
360         """
361         odir = os.path.join(self.outputdir, os.path.basename(self.pre))
362         pretest = Cmd(self.pre, outputdir=odir, timeout=self.timeout,
363                      user=self.pre_user)
364         test = Cmd(self.pathname, outputdir=self.outputdir,
365                  timeout=self.timeout, user=self.user)
366         odir = os.path.join(self.outputdir, os.path.basename(self.post))
367         posttest = Cmd(self.post, outputdir=odir, timeout=self.timeout,
368                       user=self.post_user)

370         cont = True
371         if len(pretest.pathname):
372             pretest.run(options)
373             cont = pretest.result is 'PASS'
374             pretest.log(logger, options)

376         if cont:
377             test.run(options)
378         else:
379             test.skip()

381         test.log(logger, options)

383         if len(posttest.pathname):
384             posttest.run(options)
385             posttest.log(logger, options)

388 class TestGroup(Test):
389     props = Test.props + ['tests']

391     def __init__(self, pathname, outputdir=None, timeout=None, user=None,

```



```

392         pre=None, pre_user=None, post=None, post_user=None,
393         tests=None):
394     super(TestGroup, self).__init__(pathname, outputdir, timeout, user,
395         pre, pre_user, post, post_user)
396     self.tests = tests or []

398     def __str__(self):
399         post_user = pre_user = ''
400         if len(self.pre_user):
401             pre_user = ' (as %s)' % (self.pre_user)
402         if len(self.post_user):
403             post_user = ' (as %s)' % (self.post_user)
404         return "Pathname: %s\nOutputdir: %s\nTests: %s\nTimeout: %d\n \
405             \"Pre: %s%s\nPost: %s%s\nUser: %s\n\" % \
406             (self.pathname, self.outputdir, self.tests, self.timeout,
407             self.pre, pre_user, self.post, post_user, self.user)

409     def verify(self, logger):
410         """
411         Check the pre/post scripts, user and tests in this TestGroup. Omit
412         the TestGroup entirely, or simply delete the relevant tests in the
413         group, if that's all that's required.
414         """
415         # If the pre or post scripts are relative pathnames, convert to
416         # absolute, so they stand a chance of passing verification.
417         if len(self.pre) and not os.path.isabs(self.pre):
418             self.pre = os.path.join(self.pathname, self.pre)
419         if len(self.post) and not os.path.isabs(self.post):
420             self.post = os.path.join(self.pathname, self.post)

422         auxfiles = [self.pre, self.post]
423         users = [self.pre_user, self.user, self.post_user]

425         for f in [f for f in auxfiles if len(f)]:
426             if self.pathname != os.path.dirname(f):
427                 logger.info("Warning: TestGroup '%s' not added to this run. "
428                     "Auxiliary script '%s' exists in a different "
429                     "directory." % (self.pathname, f))
430                 return False

432             if not verify_file(f):
433                 logger.info("Warning: TestGroup '%s' not added to this run. "
434                     "Auxiliary script '%s' failed verification." %
435                     (self.pathname, f))
436                 return False

438         for user in [user for user in users if len(user)]:
439             if not verify_user(user, logger):
440                 logger.info("Not adding TestGroup '%s' to this run." %
441                     self.pathname)
442                 return False

444         # If one of the tests is invalid, delete it, log it, and drive on.
445         self.tests[:] = [f for f in self.tests if
446             verify_file(os.path.join(self.pathname, f))]

448         return len(self.tests) is not 0

450     def run(self, logger, options):
451         """
452         Create Cmd instances for the pre/post scripts. If the pre script
453         doesn't pass, skip all the tests in this TestGroup. Run the post
454         script regardless.
455         """
456         odir = os.path.join(self.outputdir, os.path.basename(self.pre))
457         pretest = Cmd(self.pre, outputdir=odir, timeout=self.timeout,

```

```

458         user=self.pre_user)
459         odir = os.path.join(self.outputdir, os.path.basename(self.post))
460         posttest = Cmd(self.post, outputdir=odir, timeout=self.timeout,
461             user=self.post_user)

463         cont = True
464         if len(pretest.pathname):
465             pretest.run(options)
466             cont = pretest.result.result is 'PASS'
467             pretest.log(logger, options)

469         for fname in self.tests:
470             test = Cmd(os.path.join(self.pathname, fname),
471                 outputdir=os.path.join(self.outputdir, fname),
472                 timeout=self.timeout, user=self.user)
473             if cont:
474                 test.run(options)
475             else:
476                 test.skip()

478             test.log(logger, options)

480         if len(posttest.pathname):
481             posttest.run(options)
482             posttest.log(logger, options)

485 class TestRun(object):
486     props = ['quiet', 'outputdir']

488     def __init__(self, options):
489         self.tests = {}
490         self.testgroups = {}
491         self.starttime = time()
492         self.timestamp = datetime.now().strftime('%Y%m%dT%H%M%S')
493         self.outputdir = os.path.join(options.outputdir, self.timestamp)
494         self.logger = self.setup_logging(options)
495         self.defaults = [
496             ('outputdir', BASEDIR),
497             ('quiet', False),
498             ('timeout', 60),
499             ('user', ''),
500             ('pre', ''),
501             ('pre_user', ''),
502             ('post', ''),
503             ('post_user', '')
504         ]

506     def __str__(self):
507         s = 'TestRun:\n    outputdir: %s\n' % self.outputdir
508         s += 'TESTS:\n'
509         for key in sorted(self.tests.keys()):
510             s += '%s%s' % (self.tests[key].__str__(), '\n')
511         s += 'TESTGROUPS:\n'
512         for key in sorted(self.testgroups.keys()):
513             s += '%s%s' % (self.testgroups[key].__str__(), '\n')
514         return s

516     def adctest(self, pathname, options):
517         """
518         Create a new Test, and apply any properties that were passed in
519         from the command line. If it passes verification, add it to the
520         TestRun.
521         """
522         test = Test(pathname)
523         for prop in Test.props:

```

```

524         setattr(test, prop, getattr(options, prop))
526     if test.verify(self.logger):
527         self.tests[pathname] = test
529     def addtestgroup(self, dirname, filenames, options):
530         """
531         Create a new TestGroup, and apply any properties that were passed
532         in from the command line. If it passes verification, add it to the
533         TestRun.
534         """
535         if dirname not in self.testgroups:
536             testgroup = TestGroup(dirname)
537             for prop in Test.props:
538                 setattr(testgroup, prop, getattr(options, prop))
540         # Prevent pre/post scripts from running as regular tests
541         for f in [testgroup.pre, testgroup.post]:
542             if f in filenames:
543                 del filenames[filenames.index(f)]
545         self.testgroups[dirname] = testgroup
546         self.testgroups[dirname].tests = sorted(filenames)
548         testgroup.verify(self.logger)
550     def read(self, logger, options):
551         """
552         Read in the specified runfile, and apply the TestRun properties
553         listed in the 'DEFAULT' section to our TestRun. Then read each
554         section, and apply the appropriate properties to the Test or
555         TestGroup. Properties from individual sections override those set
556         in the 'DEFAULT' section. If the Test or TestGroup passes
557         verification, add it to the TestRun.
558         """
559         config = ConfigParser.RawConfigParser()
560         if not len(config.read(options.runfile)):
561             fail("Couldn't read config file %s" % options.runfile)
563         for opt in TestRun.props:
564             if config.has_option('DEFAULT', opt):
565                 setattr(self, opt, config.get('DEFAULT', opt))
566         self.outputdir = os.path.join(self.outputdir, self.timestamp)
568         for section in config.sections():
569             if ('arch' in config.options(section) and
570                 platform.machine() != config.get(section, 'arch')):
571                 continue
573             if 'tests' in config.options(section):
574                 testgroup = TestGroup(section)
575                 for prop in TestGroup.props:
576                     for sect in ['DEFAULT', section]:
577                         if config.has_option(sect, prop):
578                             setattr(testgroup, prop, config.get(sect, prop))
580             # Repopulate tests using eval to convert the string to a list
581             testgroup.tests = eval(config.get(section, 'tests'))
583             if testgroup.verify(logger):
584                 self.testgroups[section] = testgroup
586             elif 'autotests' in config.options(section):
587                 testgroup = TestGroup(section)
588                 for prop in TestGroup.props:
589                     for sect in ['DEFAULT', section]:

```

```

590             if config.has_option(sect, prop):
591                 setattr(testgroup, prop, config.get(sect, prop))
593         filenames = os.listdir(section)
594         # only files starting with "tst." are considered tests
595         filenames = [f for f in filenames if f.startswith("tst.")]
596         testgroup.tests = sorted(filenames)
598         if testgroup.verify(logger):
599             self.testgroups[section] = testgroup
601     else:
602         test = Test(section)
603         for prop in Test.props:
604             for sect in ['DEFAULT', section]:
605                 if config.has_option(sect, prop):
606                     setattr(test, prop, config.get(sect, prop))
608         if test.verify(logger):
609             self.tests[section] = test
611     def write(self, options):
612         """
613         Create a configuration file for editing and later use. The
614         'DEFAULT' section of the config file is created from the
615         properties that were specified on the command line. Tests are
616         simply added as sections that inherit everything from the
617         'DEFAULT' section. TestGroups are the same, except they get an
618         option including all the tests to run in that directory.
619         """
621         defaults = dict([(prop, getattr(options, prop)) for prop, _ in
622                         self.defaults])
623         config = ConfigParser.RawConfigParser(defaults)
625         for test in sorted(self.tests.keys()):
626             config.add_section(test)
628         for testgroup in sorted(self.testgroups.keys()):
629             config.add_section(testgroup)
630             config.set(testgroup, 'tests', self.testgroups[testgroup].tests)
632         try:
633             with open(options.template, 'w') as f:
634                 return config.write(f)
635         except IOError:
636             fail('Could not open \'%s\' for writing.' % options.template)
638     def complete_outputdirs(self):
639         """
640         Collect all the pathnames for Tests, and TestGroups. Work
641         backwards one pathname component at a time, to create a unique
642         directory name in which to deposit test output. Tests will be able
643         to write output files directly in the newly modified outputdir.
644         TestGroups will be able to create one subdirectory per test in the
645         outputdir, and are guaranteed uniqueness because a group can only
646         contain files in one directory. Pre and post tests will create a
647         directory rooted at the outputdir of the Test or TestGroup in
648         question for their output.
649         """
650         done = False
651         components = 0
652         tmp_dict = dict(self.tests.items() + self.testgroups.items())
653         total = len(tmp_dict)
654         base = self.outputdir

```

```

656     while not done:
657         l = []
658         components -= 1
659         for testfile in tmp_dict.keys():
660             uniq = '/'.join(testfile.split('/')[components:]).rstrip('/')
661             if uniq not in l:
662                 l.append(uniq)
663                 tmp_dict[testfile].outputdir = os.path.join(base, uniq)
664             else:
665                 break
666         done = total == len(l)

668 def setup_logging(self, options):
669     """
670     Two loggers are set up here. The first is for the logfile which
671     will contain one line summarizing the test, including the test
672     name, result, and running time. This logger will also capture the
673     timestamped combined stdout and stderr of each run. The second
674     logger is optional console output, which will contain only the one
675     line summary. The loggers are initialized at two different levels
676     to facilitate segregating the output.
677     """
678     if options.dryrun is True:
679         return

681     testlogger = logging.getLogger(__name__)
682     testlogger.setLevel(logging.DEBUG)

684     if options.cmd is not 'wrconfig':
685         try:
686             old = os.umask(0)
687             os.makedirs(self.outputdir, mode=0777)
688             os.umask(old)
689         except OSError, e:
690             fail('%s' % e)
691         filename = os.path.join(self.outputdir, 'log')

693         logfile = WatchedFileHandlerClosed(filename)
694         logfile.setLevel(logging.DEBUG)
695         logfilefmt = logging.Formatter('%(message)s')
696         logfile.setFormatter(logfilefmt)
697         testlogger.addHandler(logfile)

699         cons = logging.StreamHandler()
700         cons.setLevel(logging.INFO)
701         consfmt = logging.Formatter('%(message)s')
702         cons.setFormatter(consfmt)
703         testlogger.addHandler(cons)

705     return testlogger

707 def run(self, options):
708     """
709     Walk through all the Tests and TestGroups, calling run().
710     """
711     if not options.dryrun:
712         try:
713             os.chdir(self.outputdir)
714         except OSError:
715             fail('Could not change to directory %s' % self.outputdir)
716     for test in sorted(self.tests.keys()):
717         self.tests[test].run(self.logger, options)
718     for testgroup in sorted(self.testgroups.keys()):
719         self.testgroups[testgroup].run(self.logger, options)

721 def summary(self):

```

```

722         if Result.total is 0:
723             return

725     print '\nResults Summary'
726     for key in Result.runresults.keys():
727         if Result.runresults[key] is not 0:
728             print '%s\t% 4d' % (key, Result.runresults[key])

730     m, s = divmod(time() - self.starttime, 60)
731     h, m = divmod(m, 60)
732     print '\nRunning Time:\t%02d:%02d:%02d' % (h, m, s)
733     print 'Percent passed:\t%.1f%%' % ((float(Result.runresults['PASS']) /
734                                         float(Result.total)) * 100)
735     print 'Log directory:\t%s' % self.outputdir

738 def verify_file(pathname):
739     """
740     Verify that the supplied pathname is an executable regular file.
741     """
742     if os.path.isdir(pathname) or os.path.islink(pathname):
743         return False

745     if os.path.isfile(pathname) and os.access(pathname, os.X_OK):
746         return True

748     return False

751 def verify_user(user, logger):
752     """
753     Verify that the specified user exists on this system, and can execute
754     sudo without being prompted for a password.
755     """
756     testcmd = [SUDO, '-n', '-u', user, TRUE]

758     if user in Cmd.verified_users:
759         return True

761     try:
762         _ = getpwnam(user)
763     except KeyError:
764         logger.info("Warning: user '%s' does not exist.", user)
765         return False

767     p = Popen(testcmd)
768     p.wait()
769     if p.returncode is not 0:
770         logger.info("Warning: user '%s' cannot use passwordless sudo.", user)
771         return False
772     else:
773         Cmd.verified_users.append(user)

775     return True

778 def find_tests(testrun, options):
779     """
780     For the given list of pathnames, add files as Tests. For directories,
781     if do_groups is True, add the directory as a TestGroup. If False,
782     recursively search for executable files.
783     """

785     for p in sorted(options.pathnames):
786         if os.path.isdir(p):
787             for dirname, _, filenames in os.walk(p):

```

```

788         if options.do_groups:
789             testrun.addtestgroup(dirname, filenames, options)
790         else:
791             for f in sorted(filenames):
792                 testrun.addtest(os.path.join(dirname, f), options)
793     else:
794         testrun.addtest(p, options)

797 def fail(retstr, ret=1):
798     print '%s: %s' % (argv[0], retstr)
799     exit(ret)

802 def options_cb(option, opt_str, value, parser):
803     path_options = ['runfile', 'outputdir', 'template']

805     if option.dest is 'runfile' and '-w' in parser.rargs or \
806        option.dest is 'template' and '-c' in parser.rargs:
807         fail('-c and -w are mutually exclusive.')

809     if opt_str in parser.rargs:
810         fail('%s may only be specified once.' % opt_str)

812     if option.dest is 'runfile':
813         parser.values.cmd = 'rdconfig'
814     if option.dest is 'template':
815         parser.values.cmd = 'wrconfig'

817     setattr(parser.values, option.dest, value)
818     if option.dest in path_options:
819         setattr(parser.values, option.dest, os.path.abspath(value))

822 def parse_args():
823     parser = OptionParser()
824     parser.add_option('-c', action='callback', callback=options_cb,
825                      type='string', dest='runfile', metavar='runfile',
826                      help='Specify tests to run via config file.')
827     parser.add_option('-d', action='store_true', default=False, dest='dryrun',
828                      help='Dry run. Print tests, but take no other action.')
829     parser.add_option('-g', action='store_true', default=False,
830                      dest='do_groups', help='Make directories TestGroups.')
831     parser.add_option('-o', action='callback', callback=options_cb,
832                      default=BASEDIR, dest='outputdir', type='string',
833                      metavar='outputdir', help='Specify an output directory.')
834     parser.add_option('-p', action='callback', callback=options_cb,
835                      default='', dest='pre', metavar='script',
836                      type='string', help='Specify a pre script.')
837     parser.add_option('-P', action='callback', callback=options_cb,
838                      default='', dest='post', metavar='script',
839                      type='string', help='Specify a post script.')
840     parser.add_option('-q', action='store_true', default=False, dest='quiet',
841                      help='Silence on the console during a test run.')
842     parser.add_option('-t', action='callback', callback=options_cb, default=60,
843                      dest='timeout', metavar='seconds', type='int',
844                      help='Timeout (in seconds) for an individual test.')
845     parser.add_option('-u', action='callback', callback=options_cb,
846                      default='', dest='user', metavar='user', type='string',
847                      help='Specify a different user name to run as.')
848     parser.add_option('-w', action='callback', callback=options_cb,
849                      default=None, dest='template', metavar='template',
850                      type='string', help='Create a new config file.')
851     parser.add_option('-x', action='callback', callback=options_cb, default='',
852                      dest='pre_user', metavar='pre_user', type='string',
853                      help='Specify a user to execute the pre script.')

```

```

854     parser.add_option('-X', action='callback', callback=options_cb, default='',
855                      dest='post_user', metavar='post_user', type='string',
856                      help='Specify a user to execute the post script.')
857     (options, pathnames) = parser.parse_args()

859     if not options.runfile and not options.template:
860         options.cmd = 'runtests'

862     if options.runfile and len(pathnames):
863         fail('Extraneous arguments.')

865     options.pathnames = [os.path.abspath(path) for path in pathnames]

867     return options

870 def main():
871     options = parse_args()
872     testrun = TestRun(options)

874     if options.cmd is 'runtests':
875         find_tests(testrun, options)
876     elif options.cmd is 'rdconfig':
877         testrun.read(testrun.logger, options)
878     elif options.cmd is 'wrconfig':
879         find_tests(testrun, options)
880         testrun.write(options)
881         exit(0)
882     else:
883         fail('Unknown command specified')

885     testrun.complete_outputdirs()
886     testrun.run(options)
887     testrun.summary()
888     exit(0)

891 if __name__ == '__main__':
892     main()

```

```

*****
29753 Thu Jun 14 17:15:58 2018
new/usr/src/uts/common/sys/proc.h
9600 LDT still not happy under KPTI
*****
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) 1988, 2010, Oracle and/or its affiliates. All rights reserved.
24  * Copyright 2018 Joyent, Inc.
25  * Copyright 2017 Joyent, Inc.
26 */

27 /*      Copyright (c) 1984, 1986, 1987, 1988, 1989 AT&T */
28 /*      All Rights Reserved */

30 #ifndef _SYS_PROC_H
31 #define _SYS_PROC_H

33 #include <sys/time.h>
34 #include <sys/thread.h>
35 #include <sys/cred.h>
36 #include <sys/user.h>
37 #include <sys/watchpoint.h>
38 #include <sys/timer.h>
39 #if defined(__x86)
40 #include <sys/tss.h>
41 #include <sys/segments.h>
42 #endif
43 #include <sys/utrap.h>
44 #include <sys/model.h>
45 #include <sys/refstr.h>
46 #include <sys/avl.h>
47 #include <sys/rctl.h>
48 #include <sys/list.h>
49 #include <sys/avl.h>
50 #include <sys/door_impl.h>
51 #include <sys/signalfd.h>
52 #include <sys/secflags.h>

54 #ifdef __cplusplus
55 extern "C" {
56 #endif

58 /*
59  * Profile arguments.
60 */

```

```

61 struct prof {
62     void                *pr_base;        /* buffer base */
63     uintptr_t           pr_off;          /* pc offset */
64     size_t               pr_size;        /* buffer size */
65     uint32_t            pr_scale;        /* pc scaling */
66     long                 pr_samples;     /* sample count */
67 };
unchanged_portion_omitted

127 struct pool;
128 struct task;
129 struct zone;
130 struct brand;
131 struct corectl_path;
132 struct corectl_content;

134 /*
135  * One structure allocated per active process. It contains all
136  * data needed about the process while the process may be swapped
137  * out. Other per-process data (user.h) is also inside the proc structure.
138  * Lightweight-process data (lwp.h) and the kernel stack may be swapped out.
139  */
140 typedef struct proc {
141     /*
142      * Fields requiring no explicit locking
143      */
144     struct vnode *p_exec;                /* pointer to a.out vnode */
145     struct as *p_as;                     /* process address space pointer */
146     struct plock *p_lockp;               /* ptr to proc struct's mutex lock */
147     kmutex_t p_crlock;                   /* lock for p_cred */
148     struct cred *p_cred;                 /* process credentials */
149     /*
150      * Fields protected by pidlock
151      */
152     int p_swapcnt;                       /* number of swapped out lwps */
153     char p_stat;                          /* status of process */
154     char p_wcode;                         /* current wait code */
155     ushort_t p_pidflag;                   /* flags protected only by pidlock */
156     int p_wdata;                          /* current wait return value */
157     pid_t p_ppid;                         /* process id of parent */
158     struct proc *p_link;                  /* forward link */
159     struct proc *p_parent;                /* ptr to parent process */
160     struct proc *p_child;                 /* ptr to first child process */
161     struct proc *p_sibling;               /* ptr to next sibling proc on chain */
162     struct proc *p_psibling;              /* ptr to prev sibling proc on chain */
163     struct proc *p_sibling_ns;           /* prt to siblings with new state */
164     struct proc *p_child_ns;             /* prt to children with new state */
165     struct proc *p_next;                  /* active chain link next */
166     struct proc *p_prev;                  /* active chain link prev */
167     struct proc *p_nextofkin;            /* gets accounting info at exit */
168     struct proc *p_orphan;
169     struct proc *p_nextorph;
170     struct proc *p_pglink;                /* process group hash chain link next */
171     struct proc *p_ppglink;              /* process group hash chain link prev */
172     struct sess *p_sessp;                 /* session information */
173     struct pid *p_pidp;                   /* process ID info */
174     struct pid *p_pgidp;                  /* process group ID info */
175     /*
176      * Fields protected by p_lock
177      */
178     kcondvar_t p_cv;                      /* proc struct's condition variable */
179     kcondvar_t p_flag_cv;
180     kcondvar_t p_lwpexit;                 /* waiting for some lwp to exit */
181     kcondvar_t p_holdlwps;                /* process is waiting for its lwps */
182     /* to be held. */
183     uint_t p_proc_flag;                   /* /proc-related flags */

```

```

184     uint_t    p_flag;                /* protected while set. */
185                                     /* flags defined below */
186     clock_t  p_utime;                /* user time, this process */
187     clock_t  p_stime;                /* system time, this process */
188     clock_t  p_cutime;               /* sum of children's user time */
189     clock_t  p_cstime;               /* sum of children's system time */
190     avl_tree_t *p_segacct;           /* System V shared segment list */
191     avl_tree_t *p_semacct;           /* System V semaphore undo list */
192     caddr_t  p_bssbase;               /* base addr of last bss below heap */
193     caddr_t  p_brkbase;               /* base addr of heap */
194     size_t   p_brksize;               /* heap size in bytes */
195     uint_t   p_brkpagesz;            /* preferred heap max page size code */
196     /*
197     * Per process signal stuff.
198     */
199     k_sigset_t p_sig;                 /* signals pending to this process */
200     k_sigset_t p_extsig;               /* signals sent from another contract */
201     k_sigset_t p_ignore;               /* ignore when generated */
202     k_sigset_t p_siginfo;              /* gets signal info with signal */
203     void *p_sigfd;                    /* signalfd support state */
204     struct sigqueue *p_sigqueue;       /* queued siginfo structures */
205     struct sigqhdr *p_sigqhdr;         /* hdr to sigqueue structure pool */
206     struct sigqhdr *p_sighdr;         /* hdr to signotify structure pool */
207     uchar_t  p_stopsig;                /* jobcontrol stop signal */

209     /*
210     * Special per-process flag when set will fix misaligned memory
211     * references.
212     */
213     char     p_fixalignment;

215     /*
216     * Per process lwp and kernel thread stuff
217     */
218     id_t     p_lwpid;                  /* most recently allocated lwpid */
219     int      p_lwpcnt;                 /* number of lwps in this process */
220     int      p_lwprcnt;                 /* number of not stopped lwps */
221     int      p_lwpdaemon;              /* number of TP_DAEMON lwps */
222     int      p_lwpwait;                 /* number of lwps in lwp_wait() */
223     int      p_lwpdwait;                /* number of daemons in lwp_wait() */
224     int      p_zombcnt;                 /* number of zombie lwps */
225     kthread_t *p_tlist;                 /* circular list of threads */
226     lwpdir_t *p_lwpdir;                 /* thread (lwp) directory */
227     lwpdir_t *p_lwpfree;                 /* p_lwpdir free list */
228     tidhash_t *p_tidhash;               /* tid (lwpid) lookup hash table */
229     uint_t   p_lwpdir_sz;                /* number of p_lwpdir[] entries */
230     uint_t   p_tidhash_sz;              /* number of p_tidhash[] entries */
231     ret_tidhash_t *p_ret_tidhash;       /* retired tidhash hash tables */
232     uint64_t p_lgrpset;                  /* unprotected hint of set of lgrps */
233     /* on which process has threads */
234     volatile lgrp_id_t p_tl_lgrp;        /* main's thread lgroup id */
235     volatile lgrp_id_t p_tr_lgrp;        /* text replica's lgroup id */
236 #if defined(LP64)
237     uintptr_t p_lgrpres2;                /* reserved for lgrp migration */
238 #endif

239     /*
240     * /proc (process filesystem) debugger interface stuff.
241     */
242     k_sigset_t p_sigmask;                /* mask of traced signals (/proc) */
243     k_fltset_t p_fltmask;                /* mask of traced faults (/proc) */
244     struct vnode *p_trace;                /* pointer to primary /proc vnode */
245     struct vnode *p_plist;                /* list of /proc vnodes for process */
246     kthread_t *p_agenttp;                 /* thread ptr for /proc agent lwp */
247     avl_tree_t p_warea;                  /* list of watched areas */
248     avl_tree_t p_wpage;                  /* remembered watched pages (vfork) */
249     watched_page_t *p_wprot;              /* pages that need to have prot set */

```

```

250     int      p_mapcnt;                  /* number of active pr_mappage()s */
251     kmutex_t p_maplock;                 /* lock for pr_mappage() */
252     struct proc *p_rlink;                /* linked list for server */
253     kcondvar_t p_srwchan_cv;
254     size_t   p_stksize;                  /* process stack size in bytes */
255     uint_t   p_stkpagesz;                /* preferred stack max page size code */

257     /*
258     * Microstate accounting, resource usage, and real-time profiling
259     */
260     hrtime_t p_mstart;                    /* hi-res process start time */
261     hrtime_t p_mterm;                     /* hi-res process termination time */
262     hrtime_t p_mlreal;                    /* elapsed time sum over defunct lwps */
263     hrtime_t p_acct[NMSTATES];            /* microstate sum over defunct lwps */
264     hrtime_t p_cacct[NMSTATES];          /* microstate sum over child procs */
265     struct lrusage p_ru;                  /* lrusage sum over defunct lwps */
266     struct lrusage p_cru;                 /* lrusage sum over child procs */
267     struct itimerval p_rprof_timer;        /* ITIMER_REALPROF interval timer */
268     uintptr_t p_rprof_cyclic;            /* ITIMER_REALPROF cyclic */
269     uint_t   p_defunct;                  /* number of defunct lwps */
270     /*
271     * profiling. A lock is used in the event of multiple lwp's
272     * using the same profiling base/size.
273     */
274     kmutex_t p_plock;                    /* protects user profile arguments */
275     struct prof p_prof;                   /* profile arguments */

277     /*
278     * Doors.
279     */
280     door_pool_t p_server_threads;         /* common thread pool */
281     struct door_node *p_door_list;        /* active doors */
282     struct door_node *p_unref_list;       /* p_unref_cv;
283     kcondvar_t p_unref_cv;                /* unref thread created */
284     char

286     /*
287     * Kernel probes
288     */
289     uchar_t p_tnf_flags;

291     /*
292     * Solaris Audit
293     */
294     struct p_audit_data *p_audit_data;    /* per process audit structure */

296     pctxop_t *p_pctx;

298 #if defined(__x86)
299     /*
300     * LDT support.
301     */
302     kmutex_t p_ldtlock;                  /* protects the following fields */
303     user_desc_t *p_ldt;                   /* Pointer to private LDT */
304     uint64_t p_unused1;                   /* no longer used */
305     uint64_t p_unused2;                   /* no longer used */
306     system_desc_t *p_ldt_desc;            /* segment descriptor for private LDT */
307     ushort_t p_ldtlimit;                  /* highest selector used */
308 #endif

309     size_t p_swrss;                       /* resident set size before last swap */
310     struct aio *p_aio;                    /* pointer to async I/O struct */
311     struct itimer **p_itimer;              /* interval timers */
312     timeout_id_t p_alarmid;                /* alarm's timeout id */
313     caddr_t p_usrstack;                   /* top of the process stack */
314     uint_t p_stkprot;                      /* stack memory protection */
315     uint_t p_datprot;                      /* data memory protection */

```

```

315     model_t      p_model;          /* data model determined at exec time */
316     struct lwpchan_data *p_lcp;    /* lwpchan cache */
317     kmutex_t     p_lcp_lock;      /* protects assignments to p_lcp */
318     utrap_handler_t *p_utrap;     /* pointer to user trap handlers */
319     struct corectl_path *p_corefile; /* pattern for core file */
320     struct task  *p_task;          /* our containing task */
321     struct proc  *p_taskprev;     /* ptr to previous process in task */
322     struct proc  *p_tasknext;     /* ptr to next process in task */
323     kmutex_t     p_sc_lock;       /* protects p_pagep */
324     struct sc_page_ctl *p_pagep;  /* list of process's shared pages */
325     struct rctl_set *p_rctls;     /* resource controls for this process */
326     rlim64_t     p_stk_ctl;       /* currently enforced stack size */
327     rlim64_t     p_fsz_ctl;       /* currently enforced file size */
328     rlim64_t     p_vmем_ctl;      /* currently enforced addr-space size */
329     rlim64_t     p_fno_ctl;       /* currently enforced file-desc limit */
330     pid_t        p_ancpid;        /* ancestor pid, used by exactt */
331     struct itimerval p_realtimer; /* real interval timer */
332     timeout_id_t  p_itimerid;     /* real interval timer's timeout id */
333     struct corectl_content *p_content; /* content of core file */

335     avl_tree_t    p_ct_held;       /* held contracts */
336     struct ct_equeue **p_ct_equeue; /* process-type event queues */

338     struct cont_process *p_ct_process; /* process contract */
339     list_node_t        p_ct_member;   /* process contract membership */
340     sigqueue_t        *p_killsgp;    /* sigqueue pointer for SIGKILL */

342     int                p_dtrace_probes; /* are there probes for this proc? */
343     uint64_t           p_dtrace_count; /* number of DTrace tracepoints */
344                                     /* (protected by P_PR_LOCK) */
345     void                *p_dtrace_helpers; /* DTrace helpers, if any */
346     struct pool         *p_pool;        /* pointer to containing pool */
347     kcondvar_t          p_poolcv;      /* synchronization with pools */
348     uint_t              p_poolcnt;     /* # threads inside pool barrier */
349     uint_t              p_poolflag;    /* pool-related flags (see below) */
350     uintptr_t           p_portcnt;     /* event ports counter */
351     struct zone         *p_zone;       /* zone in which process lives */
352     struct vnode        *p_execdir;    /* directory that p_exec came from */
353     struct brand        *p_brand;      /* process's brand */
354     void                *p_brand_data; /* per-process brand state */
355     psecflags_t         p_secflags;    /* per-process security flags */

357     /* additional lock to protect p_sespp (but not its contents) */
358     kmutex_t p_spllock;
359     rctl_qty_t p_locked_mem; /* locked memory charged to proc */
360                                     /* protected by p_lock */
361     rctl_qty_t p_crypto_mem; /* /dev/crypto memory charged to proc */
362                                     /* protected by p_lock */
363     clock_t p_tptime; /* buffered task time */

365     /*
366     * The user structure
367     */
368     struct user p_user; /* (see sys/user.h) */
369 } proc_t;

```

unchanged portion omitted

```

*****
8657 Thu Jun 14 17:15:58 2018
new/usr/src/uts/i86pc/ml/offsets.in
9600 LDT still not happy under KPTI
*****
1 \
2 \ Copyright (c) 2004, 2010, Oracle and/or its affiliates. All rights reserved.
3 \ Copyright 2012 Garrett D'Amore <garrett@damore.org>. All rights reserved.
4 \ Copyright 2018 Joyent, Inc.
5 \
6 \ CDDL HEADER START
7 \
8 \ The contents of this file are subject to the terms of the
9 \ Common Development and Distribution License (the "License").
10 \ You may not use this file except in compliance with the License.
11 \
12 \ You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
13 \ or http://www.opensolaris.org/os/licensing.
14 \ See the License for the specific language governing permissions
15 \ and limitations under the License.
16 \
17 \ When distributing Covered Code, include this CDDL HEADER in each
18 \ file and include the License file at usr/src/OPENSOLARIS.LICENSE.
19 \ If applicable, add the following below this CDDL HEADER, with the
20 \ fields enclosed by brackets "[]" replaced with your own identifying
21 \ information: Portions Copyright [yyyy] [name of copyright owner]
22 \
23 \ CDDL HEADER END
24 \

27 \
28 \ offsets.in: input file to produce assym.h using the ctfstabs program
29 \

31 #ifndef _GENASSYM
32 #define _GENASSYM
33 #endif

35 #define SIZES 1

37 #include <sys/types.h>
38 #include <sys/bootsvcs.h>
39 #include <sys/system.h>
40 #include <sys/sysinfo.h>
41 #include <sys/user.h>
42 #include <sys/thread.h>
43 #include <sys/proc.h>
44 #include <sys/cpuvar.h>
45 #include <sys/tss.h>
46 #include <sys/privregs.h>
47 #include <sys/segments.h>
48 #include <sys/devops.h>
49 #include <sys/ddi_impldefs.h>
50 #include <vm/as.h>
51 #include <sys/avintr.h>
52 #include <sys/pic.h>
53 #include <sys/rm_platter.h>
54 #include <sys/stream.h>
55 #include <sys/strsubr.h>
56 #include <sys/sunddi.h>
57 #include <sys/modctl.h>
58 #include <sys/traptrace.h>
59 #include <sys/traptrap.h>
60 #include <sys/lgrp.h>
61 #include <sys/dtrace.h>

```

```

62 #include <sys/brand.h>
63 #include <sys/fastboot.h>
64 #include <sys/cpr_wakecode.h>
65 #include <sys/comm_page.h>

67 proc          PROC_SIZE
68     p_link
69     p_next
70     p_child
71     p_sibling
72     p_sig
73     p_flag
74     p_tlist
75     p_as
76     p_lockp
77     p_user
78     p_ldt
79     p_ldt_desc
78     p_model
79     p_pctx
80     p_agenttp
81     p_zone
82     p_brand
83     p_brand_data

85 _kthread      THREAD_SIZE
86     t_pcb
87     t_lock
88     t_lockstat
89     t_lockp
90     t_lock_flush
91     t_kpri_req
92     t_oldspl
93     t_pri
94     t_pil
95     t_lwp
96     t_procp
97     t_link
98     t_state
99     t_mstate
100    t_preempt_lk
101    t_stk
102    t_swap
103    t_lwpchan.lc_wchan
104    t_flag
105    t_ctx
106    t_lofault
107    t_onfault
108    t_ontrap
109    t_cpu
110    t_lpl
111    t_bound_cpu
112    t_intr
113    t_forw
114    t_back
115    t_sig
116    t_tid
117    t_pre_sys
118    t_preempt
119    t_proc_flag
120    t_startpc
121    t_sysnum
122    t_intr_start
123    _tu._ts._t_astflag
124    _tu._ts._t_post_sys
125    _tu._t_post_sys_ast

T_LABEL
T_STACK
T_WCHAN
T_FLAGS
T_ASTFLAG
T_POST_SYS
T_POST_SYS_AST

```



```

126     t_copyops
127 #ifdef  __amd64
128     t_useracc
129 #endif

131 ctxop
132     save_op          CTXOP_SAVE

134 as
135     a_hat

137 user    USIZEBYTES
138     u_comm
139     u_signal

141 _label_t
142     val      LABEL_VAL

144 #define    LABEL_PC      LABEL_VAL
145 #define    LABEL_SP      _CONST(LABEL_VAL + LABEL_VAL_INCR)
146 #define    T_PC          _CONST(T_LABEL + LABEL_PC)
147 #define    T_SP          _CONST(T_LABEL + LABEL_SP)

149 _klwp
150     lwp_thread
151     lwp_procp
152     lwp_brand
153     lwp_eosys
154     lwp_regs
155     lwp_arg
156     lwp_ap
157     lwp_cursig
158     lwp_state
159     lwp_mstate.ms_acct    LWP_MS_ACCT
160     lwp_mstate.ms_prev   LWP_MS_PREV
161     lwp_mstate.ms_start  LWP_MS_START
162     lwp_mstate.ms_state_start LWP_MS_STATE_START
163     lwp_pcb
164     lwp_ru.sysc          LWP_RU_SYSC

166 #define    LWP_ACCT_USER    _CONST(LWP_MS_ACCT + _MUL(LMS_USER, LWP_MS_ACCT_)
167 #define    LWP_ACCT_SYSTEM _CONST(LWP_MS_ACCT + _MUL(LMS_SYSTEM, LWP_MS_ACC

169 fpu_ctx
170     fpu_regs          FPU_CTX_FPU_REGS
171     fpu_flags         FPU_CTX_FPU_FLAGS
172     fpu_xsave_mask   FPU_CTX_FPU_XSAVE_MASK

174 fxsave_state    FXSAVE_STATE_SIZE
175     fx_fsw        FXSAVE_STATE_FSW
176     fx_mxcsr_mask FXSAVE_STATE_MXCSR_MASK

179 autovec        AUTOVECSIZE
180     av_vector
181     av_intarg1
182     av_intarg2
183     av_ticksp
184     av_link
185     av_prilevel
186     av_dip

188 av_head
189     avh_link
190     avh_hi_pri
191     avh_lo_pri

```

```

193 cpu
194     cpu_id
195     cpu_flags
196     cpu_self
197     cpu_thread
198     cpu_thread_lock
199     cpu_kprunrun
200     cpu_lwp
201     cpu_fpowner
202     cpu_idle_thread
203     cpu_intr_thread
204     cpu_intr_actv
205     cpu_base_spl
206     cpu_intr_stack
207     cpu_stats.sys.cpumigrate    CPU_STATS_SYS_CPUMIGRATE
208     cpu_stats.sys.intr          CPU_STATS_SYS_INTR
209     cpu_stats.sys.intrblk       CPU_STATS_SYS_INTRBLK
210     cpu_stats.sys.syscall       CPU_STATS_SYS_SYSCALL
211     cpu_profile_pc
212     cpu_profile_upc
213     cpu_profile_pil
214     cpu_ftrace.ftd_state        CPU_FTRACE_STATE
215     cpu_mstate
216     cpu_intracct

218 #define    CPU_INTR_ACTV_REF    _CONST(CPU_INTR_ACTV + 2)

220 cpu
221     cpu_m.pil_high_start    CPU_PIL_HIGH_START
222     cpu_m.intrstat          CPU_INTRSTAT
223     cpu_m.mcpu_current_hat  CPU_CURRENT_HAT
224     cpu_m.mcpu_gdt          CPU_GDT
225     cpu_m.mcpu_idt          CPU_IDT
226     cpu_m.mcpu_tss          CPU_TSS
227     cpu_m.mcpu_softinfo     CPU_SOFTINFO
228     cpu_m.mcpu_pri          CPU_PRI
229 #if defined(__xpv)
230     cpu_m.mcpu_vcpu_info    CPU_VCPU_INFO
231 #endif

233 cpu
234     cpu_m.mcpu_kpti.kf_kernel_cr3    CPU_KPTI_KCR3
235     cpu_m.mcpu_kpti.kf_user_cr3      CPU_KPTI_UCR3
236     cpu_m.mcpu_kpti.kf_tr_rsp        CPU_KPTI_TR_RSP
237     cpu_m.mcpu_kpti.kf_tr_cr3        CPU_KPTI_TR_CR3
238     cpu_m.mcpu_kpti.kf_r13           CPU_KPTI_R13
239     cpu_m.mcpu_kpti.kf_r14           CPU_KPTI_R14
240     cpu_m.mcpu_kpti.kf_tr_ret_rsp    CPU_KPTI_RET_RSP

242     cpu_m.mcpu_kpti.kf_ss             CPU_KPTI_SS
243     cpu_m.mcpu_kpti.kf_rsp            CPU_KPTI_RSP
244     cpu_m.mcpu_kpti.kf_rflags         CPU_KPTI_RFLAGS
245     cpu_m.mcpu_kpti.kf_cs             CPU_KPTI_CS
246     cpu_m.mcpu_kpti.kf_rip            CPU_KPTI_RIP
247     cpu_m.mcpu_kpti.kf_err            CPU_KPTI_ERR

249     cpu_m.mcpu_pad2                   CPU_KPTI_START
250     cpu_m.mcpu_pad3                   CPU_KPTI_END

252     cpu_m.mcpu_kpti_dbg                CPU_KPTI_DBG

254 kpti_frame
255     kf_r14          KPTI_R14
256     kf_r13          KPTI_R13
257     kf_err          KPTI_ERR

```

## new/usr/src/uts/i86pc/ml/offsets.in

5

```

258      kf_rip      KPTI_RIP
259      kf_cs       KPTI_CS
260      kf_rflags   KPTI_RFLAGS
261      kf_rsp      KPTI_RSP
262      kf_ss       KPTI_SS

264      kf_tr_rsp   KPTI_TOP

266      kf_kernel_cr3 KPTI_KCR3
267      kf_user_cr3  KPTI_UCR3
268      kf_tr_ret_rsp KPTI_RET_RSP
269      kf_tr_cr3    KPTI_TR_CR3

271      kf_tr_flag   KPTI_FLAG

273 standard_pic
274      c_curmask
275      c_iplmask

277 ddi_dma_impl
278      dmai_rflags
279      dmai_rdlip

281 dev_info
282      devi_ops      DEVI_DEV_OPS
283      devi_bus_ctl
284      devi_bus_dma_ctl
285      devi_bus_dma_allochdl
286      devi_bus_dma_freehdl
287      devi_bus_dma_bindhdl
288      devi_bus_dma_unbindhdl
289      devi_bus_dma_flush
290      devi_bus_dma_win

292 dev_ops
293      devo_bus_ops      DEVI_BUS_OPS

295 bus_ops
296      bus_ctl      OPS_CTL
297      bus_dma_map  OPS_MAP
298      bus_dma_ctl  OPS_MCTL
299      bus_dma_allochdl
300      bus_dma_freehdl
301      bus_dma_bindhdl
302      bus_dma_unbindhdl
303      bus_dma_flush
304      bus_dma_win  OPS_WIN

306 sysent  SYSENT_SIZE      SYSENT_SIZE_SHIFT
307      sy_callc
308      sy_flags
309      sy_narg

311 stdata
312      sd_lock

314 queue
315      q_flag
316      q_next
317      q_stream
318      q_syncq
319      q_qinfo

321 qinit
322      qi_putp

```

## new/usr/src/uts/i86pc/ml/offsets.in

6

```

324 syncq
325      sq_flags
326      sq_count
327      sq_lock
328      sq_wait

330 rm_platter
331      rm_idt_lim      IDTROFF
332      rm_gdt_lim      GDTRROFF
333      rm_pdbr        CR3OFF
334      rm_cpu         CPUNOFF
335      rm_cr4         CR4OFF
336      rm_cpu_halt_code CPUHALTCODEOFF
337      rm_cpu_halted  CPUHALTEDOFF

339 ddi_acc_impl
340      ahi_acc_attr   ACC_ATTR
341      ahi_get8      ACC_GETB
342      ahi_get16     ACC_GETW
343      ahi_get32     ACC_GETL
344      ahi_get64     ACC_GETLL
345      ahi_put8      ACC_PUTB
346      ahi_put16     ACC_PUTW
347      ahi_put32     ACC_PUTL
348      ahi_put64     ACC_PUTLL
349      ahi_rep_get8  ACC_REP_GETB
350      ahi_rep_get16 ACC_REP_GETW
351      ahi_rep_get32 ACC_REP_GETL
352      ahi_rep_get64 ACC_REP_GETLL
353      ahi_rep_put8  ACC_REP_PUTB
354      ahi_rep_put16 ACC_REP_PUTW
355      ahi_rep_put32 ACC_REP_PUTL
356      ahi_rep_put64 ACC_REP_PUTLL

358 on_trap_data
359      ot_prot
360      ot_trap
361      ot_trampoline
362      ot_jmpbuf
363      ot_prev
364      ot_handle
365      ot_padi

367 trap_trace_ctl_t      __TRAPTR_SIZE TRAPTR_SIZE_SHIFT
368      ttc_next          TRAPTR_NEXT
369      ttc_first         TRAPTR_FIRST
370      ttc_limit         TRAPTR_LIMIT

372 trap_trace_rec_t      TRAP_ENT_SIZE
373      ttr_cr2
374      ttr_info.idt_entry.vector      TTR_VECTOR
375      ttr_info.idt_entry.ipl        TTR_IPL
376      ttr_info.idt_entry.spl        TTR_SPL
377      ttr_info.idt_entry.pri        TTR_PRI
378      ttr_info.gate_entry.sysnum    TTR_SYSNUM
379      ttr_marker
380      ttr_stamp
381      ttr_curthread
382      ttr_sdepth
383      ttr_stack

385 lgrp_ld
386      lpl_lgrp_id

388 dtrace_id_t      DTRACE_IDSIZE

```

## new/usr/src/uts/i86pc/ml/offsets.in

7

```

390 cpu_core      CPU_CORE_SIZE  CPU_CORE_SHIFT
391   cpuc_dtrace_flags
392   cpuc_dtrace_illval

394 timespec      TIMESPEC_SIZE

396 gate_desc     GATE_DESC_SIZE

398 desctbr_t     DESC_TBR_SIZE
399   dtr_limit
400   dtr_base

402 mod_stub_info  MODS_SIZE
403   mods_func_adr  MODS_INSTFCN
404   mods_errfcn   MODS_RETFCN
405   mods_flag

407 \#define      TRAP_TSIZE      _MUL(TRAP_ENT_SIZE, TRAPTR_NENT)

409 copyops
410   cp_copyin
411   cp_xcopyin
412   cp_copyout
413   cp_xcopyout
414   cp_copyinstr
415   cp_copyoutstr
416   cp_fuword8
417   cp_fuword16
418   cp_fuword32
419   cp_fuword64
420   cp_suword8
421   cp_suword16
422   cp_suword32
423   cp_suword64
424   cp_physio

426 brand
427   b_machops

429 brand_proc_data_t
430   spd_handler

432 fastboot_file_t
433   fb_va
434   fb_pte_list_va
435   fb_pte_list_pa
436   fb_dest_pa
437   fb_size
438   fb_next_pa
439   fb_sections
440   fb_sectcnt

442 fastboot_section_t
443   fb_sec_offset
444   fb_sec_paddr
445   fb_sec_size
446   fb_sec_bss_size

448 fastboot_info_t
449   fi_files
450   fi_has_pae
451   fi_pagetable_va
452   fi_pagetable_pa
453   fi_last_table_pa
454   fi_new_mbi_pa
455   fi_valid

```

## new/usr/src/uts/i86pc/ml/offsets.in

8

```

457 zone
458   zone_brand_data

460 wc_cpu  WC_CPU_SIZE
461   wc_retaddr
462   wc_virtaddr
463   wc_cr0
464   wc_cr3
465   wc_cr4
466   wc_cr8
467   wc_fs
468   wc_fsbase
469   wc_gs
470   wc_gsbase
471   wc_kgsbase
472   wc_r8
473   wc_r9
474   wc_r10
475   wc_r11
476   wc_r12
477   wc_r13
478   wc_r14
479   wc_r15
480   wc_rax
481   wc_rbp
482   wc_rbx
483   wc_rcx
484   wc_rdi
485   wc_rdx
486   wc_rsi
487   wc_rsp
488   wc_gdt_limit  WC_GDT
489   wc_gdt_base
490   wc_idt_limit  WC_IDT
491   wc_idt_base
492   wc_tr
493   wc_ldt
494   wc_eflags
495   wc_ebx
496   wc_edi
497   wc_esi
498   wc_ebp
499   wc_esp
500   wc_esp
501   wc_ss
502   wc_cs
503   wc_ds
504   wc_es
505   wc_cpu_id
506   wc_saved_stack

508 wc_wakecode
509   wc_cpu

511 comm_page_s  COMM_PAGE_S_SIZE

```

```

*****
14793 Thu Jun 14 17:15:59 2018
new/usr/src/uts/i86pc/os/mlsetup.c
9600 LDT still not happy under KPTI
*****
_____unchanged_portion_omitted_____

98 /*
99  * Setup routine called right before main(). Interposing this function
100 * before main() allows us to call it in a machine-independent fashion.
101 */
102 void
103 mlsetup(struct regs *rp)
104 {
105     u_longlong_t prop_value;
106     extern struct classfuncs sys_classfuncs;
107     extern disp_t cpu0_disp;
108     extern char t0stack[];
109     extern int post_fastreboot;
110     extern uint64_t plat_dr_options;

112     ASSERT_STACK_ALIGNED();

114     /*
115      * initialize cpu_self
116      */
117     cpu[0]->cpu_self = cpu[0];

119 #if defined(__xpv)
120     /*
121      * Point at the hypervisor's virtual cpu structure
122      */
123     cpu[0]->cpu_m.mcpu_vcpu_info = &HYPERVISOR_shared_info->vcpu_info[0];
124 #endif

126     /*
127      * check if we've got special bits to clear or set
128      * when checking cpu features
129      */

131     if (bootprop_getval("cpuid_feature_ecx_include", &prop_value) != 0)
132         cpuid_feature_ecx_include = 0;
133     else
134         cpuid_feature_ecx_include = (uint32_t)prop_value;

136     if (bootprop_getval("cpuid_feature_ecx_exclude", &prop_value) != 0)
137         cpuid_feature_ecx_exclude = 0;
138     else
139         cpuid_feature_ecx_exclude = (uint32_t)prop_value;

141     if (bootprop_getval("cpuid_feature_edx_include", &prop_value) != 0)
142         cpuid_feature_edx_include = 0;
143     else
144         cpuid_feature_edx_include = (uint32_t)prop_value;

146     if (bootprop_getval("cpuid_feature_edx_exclude", &prop_value) != 0)
147         cpuid_feature_edx_exclude = 0;
148     else
149         cpuid_feature_edx_exclude = (uint32_t)prop_value;

151 #if !defined(__xpv)
152     /*
153      * Check to see if KPTI has been explicitly enabled or disabled.
154      * We have to check this before init_desctbls().
155      */

```

```

156     if (bootprop_getval("kpti", &prop_value) == 0) {
157         kpti_enable = (uint64_t)(prop_value == 1);
158         prom_printf("unix: forcing kpti to %s due to boot argument\n",
159                    (kpti_enable == 1) ? "ON" : "OFF");
160     } else {
161         kpti_enable = 1;
162     }

164     if (bootprop_getval("pcid", &prop_value) == 0 && prop_value == 0) {
165         prom_printf("unix: forcing pcid to OFF due to boot argument\n");
166         x86_use_pcid = 0;
167     } else if (kpti_enable != 1) {
168         x86_use_pcid = 0;
169     }
170 #endif

172     /*
173      * Initialize idt0, gdt0, ldt0_default, ktss0 and dftss.
174      */
175     init_desctbls();

177     /*
178      * lgrp_init() and possibly cpuid_pass1() need PCI config
179      * space access
180      */
181 #if defined(__xpv)
182     if (DOMAIN_IS_INITDOMAIN(xen_info))
183         pci_cfgspace_init();
184 #else
185     pci_cfgspace_init();
186     /*
187      * Initialize the platform type from CPU 0 to ensure that
188      * determine_platform() is only ever called once.
189      */
190     determine_platform();
191 #endif

193     /*
194      * The first lightweight pass (pass0) through the cpuid data
195      * was done in locore before mlsetup was called. Do the next
196      * pass in C code.
197      *
198      * The x86_featureset is initialized here based on the capabilities
199      * of the boot CPU. Note that if we choose to support CPUs that have
200      * different feature sets (at which point we would almost certainly
201      * want to set the feature bits to correspond to the feature
202      * minimum) this value may be altered.
203      */
204     cpuid_pass1(cpu[0], x86_featureset);

206 #if !defined(__xpv)
207     if ((get_hwenv() & HW_XEN_HVM) != 0)
208         xen_hvm_init();

210     /*
211      * Before we do anything with the TSCs, we need to work around
212      * Intel erratum BT81. On some CPUs, warm reset does not
213      * clear the TSC. If we are on such a CPU, we will clear TSC ourselves
214      * here. Other CPUs will clear it when we boot them later, and the
215      * resulting skew will be handled by tsc_sync_master()/_slave();
216      * note that such skew already exists and has to be handled anyway.
217      *
218      * We do this only on metal. This same problem can occur with a
219      * hypervisor that does not happen to virtualise a TSC that starts from
220      * zero, regardless of CPU type; however, we do not expect hypervisors
221      * that do not virtualise TSC that way to handle writes to TSC

```

```

222     * correctly, either.
223     */
224     if (get_hwenv() == HW_NATIVE &&
225         cpuid_getvendor(CPU) == X86_VENDOR_Intel &&
226         cpuid_getfamily(CPU) == 6 &&
227         (cpuid_getmodel(CPU) == 0x2d || cpuid_getmodel(CPU) == 0x3e) &&
228         is_x86_feature(x86_featureset, X86FSET_TSC)) {
229         (void) wrmsr(REG_TSC, 0UL);
230     }
231
232     /*
233     * Patch the tsc_read routine with appropriate set of instructions,
234     * depending on the processor family and architecture, to read the
235     * time-stamp counter while ensuring no out-of-order execution.
236     * Patch it while the kernel text is still writable.
237     *
238     * Note: tsc_read is not patched for intel processors whose family
239     * is >6 and for amd whose family >f (in case they don't support rdtscp
240     * instruction, unlikely). By default tsc_read will use cpuid for
241     * serialization in such cases. The following code needs to be
242     * revisited if intel processors of family >= f retains the
243     * instruction serialization nature of mfence instruction.
244     * Note: tsc_read is not patched for x86 processors which do
245     * not support "mfence". By default tsc_read will use cpuid for
246     * serialization in such cases.
247     *
248     * The Xen hypervisor does not correctly report whether rdtscp is
249     * supported or not, so we must assume that it is not.
250     */
251     if ((get_hwenv() & HW_XEN_HVM) == 0 &&
252         is_x86_feature(x86_featureset, X86FSET_TSCP))
253         patch_tsc_read(TSC_TSCP);
254     else if (cpuid_getvendor(CPU) == X86_VENDOR_AMD &&
255             cpuid_getfamily(CPU) <= 0xf &&
256             is_x86_feature(x86_featureset, X86FSET_SSE2))
257         patch_tsc_read(TSC_RDTSC_MFENCE);
258     else if (cpuid_getvendor(CPU) == X86_VENDOR_Intel &&
259             cpuid_getfamily(CPU) <= 6 &&
260             is_x86_feature(x86_featureset, X86FSET_SSE2))
261         patch_tsc_read(TSC_RDTSC_LFENCE);
262
263 #endif /* !__xpv */
264
265 #if defined(__i386) && !defined(__xpv)
266     /*
267     * Some i386 processors do not implement the rdtsc instruction,
268     * or at least they do not implement it correctly. Patch them to
269     * return 0.
270     */
271     if (!is_x86_feature(x86_featureset, X86FSET_TSC))
272         patch_tsc_read(TSC_NONE);
273 #endif /* __i386 && !__xpv */
274
275 #if defined(__amd64) && !defined(__xpv)
276     patch_memops(cpuid_getvendor(CPU));
277 #endif /* __amd64 && !__xpv */
278
279 #if !defined(__xpv)
280     /*
281     * While we're thinking about the TSC, let's set up %cr4 so that
282     * userland can issue rdtsc, and initialize the TSC_AUX value
283     * (the cpuid) for the rdtscp instruction on appropriately
284     * capable hardware.
285     */
286     /*
287     */

```

```

288     if (is_x86_feature(x86_featureset, X86FSET_TSC))
289         setcr4(getcr4() & ~CR4_TSD);
290
291     if (is_x86_feature(x86_featureset, X86FSET_TSCP))
292         (void) wrmsr(MSR_AMD_TSCAUX, 0);
293
294     /*
295     * Let's get the other %cr4 stuff while we're here. Note, we defer
296     * enabling CR4_SMAP until startup_end(); however, that's importantly
297     * before we start other CPUs. That ensures that it will be synced out
298     * to other CPUs.
299     */
300     if (is_x86_feature(x86_featureset, X86FSET_DE))
301         setcr4(getcr4() | CR4_DE);
302
303     if (is_x86_feature(x86_featureset, X86FSET_SMEP))
304         setcr4(getcr4() | CR4_SMEP);
305 #endif /* __xpv */
306
307     /*
308     * initialize t0
309     */
310     t0.t_stk = (caddr_t)rp - MINFRAME;
311     t0.t_stkbase = t0stack;
312     t0.t_pri = maxclsyspri - 3;
313     t0.t_schedflag = TS_LOAD | TS_DONT_SWAP;
314     t0.t_procp = &p0;
315     t0.t_plockp = &p0lock.pl_lock;
316     t0.t_lwp = &lwp0;
317     t0.t_forw = &t0;
318     t0.t_back = &t0;
319     t0.t_next = &t0;
320     t0.t_prev = &t0;
321     t0.t_cpu = cpu[0];
322     t0.t_disp_queue = &cpu0_disp;
323     t0.t_bind_cpu = PBIND_NONE;
324     t0.t_bind_pset = PS_NONE;
325     t0.t_bindflag = (uchar_t)default_binding_mode;
326     t0.t_cpupart = &cp_default;
327     t0.t_clfuncs = &sys_classfuncs.thread;
328     t0.t_copyops = NULL;
329     THREAD_ONPROC(&t0, CPU);
330
331     lwp0.lwp_thread = &t0;
332     lwp0.lwp_regs = (void *)rp;
333     lwp0.lwp_procp = &p0;
334     t0.t_tid = p0.p_lwpcnt = p0.p_lwprcnt = p0.p_lwpid = 1;
335
336     p0.p_exec = NULL;
337     p0.p_stat = SRUN;
338     p0.p_flag = SSYS;
339     p0.p_tlist = &t0;
340     p0.p_stksize = 2*PAGESIZE;
341     p0.p_stkpageszc = 0;
342     p0.p_as = &kas;
343     p0.p_lockp = &p0lock;
344     p0.p_brkpageszc = 0;
345     p0.p_tl_lgrpid = LGRP_NONE;
346     p0.p_tr_lgrpid = LGRP_NONE;
347     psecflags_default(&p0.p_secflags);
348
349     sigorset(&p0.p_ignore, &ignoredefault);
350
351     CPU->cpu_thread = &t0;
352     bzero(&cpu0_disp, sizeof (disp_t));
353     CPU->cpu_disp = &cpu0_disp;

```

```

354 CPU->cpu_disp->disp_cpu = CPU;
355 CPU->cpu_dispthread = &t0;
356 CPU->cpu_idle_thread = &t0;
357 CPU->cpu_flags = CPU_READY | CPU_RUNNING | CPU_EXISTS | CPU_ENABLE;
358 CPU->cpu_dispatch_pri = t0.t_pri;

360 CPU->cpu_id = 0;

362 CPU->cpu_pri = 12;          /* initial PIL for the boot CPU */

364 /*
365  * The kernel doesn't use LDTs unless a process explicitly requests one.
366  */
367 p0.p_ldt_desc = null_sdesc;

369 /*
370  * Initialize thread/cpu microstate accounting
371  */
372 init_mstate(&t0, LMS_SYSTEM);
373 init_cpu_mstate(CPU, CMS_SYSTEM);

375 pg_cpu_bootstrap(CPU);

377 /*
378  * Now that we have taken over the GDT, IDT and have initialized
379  * active CPU list it's time to inform kmdb if present.
380  */
381 if (boothowto & RB_DEBUG)
382     kdi_idt_sync();

384 if (BOP_GETPROPLEN(bootops, "efi-systab") < 0) {
385     /*
386      * In BIOS system, explicitly set console to text mode (0x3)
387      * if this is a boot post Fast Reboot, and the console is set
388      * to CONS_SCREEN_TEXT.
389      */
390     if (post_fastreboot &&
391         boot_console_type(NULL) == CONS_SCREEN_TEXT) {
392         set_console_mode(0x3);
393     }
394 }

396 /*
397  * If requested (boot -d) drop into kmdb.
398  *
399  * This must be done after cpu_list_init() on the 64-bit kernel
400  * since taking a trap requires that we re-compute gsbased based
401  * on the cpu list.
402  */
403 if (boothowto & RB_DEBUGENTER)
404     kmdb_enter();

406 cpu_vm_data_init(CPU);

408 rp->r_fp = 0;  /* terminate kernel stack traces! */

410 prom_init("kernel", (void *)NULL);

412 /* User-set option overrides firmware value. */
413 if (bootprop_getval(PLAT_DR_OPTIONS_NAME, &prop_value) == 0) {
414     plat_dr_options = (uint64_t)prop_value;

```

```

415     }
416 #if defined(__xpv)
417     /* No support of DR operations on xpv */
418     plat_dr_options = 0;
419 #else /* __xpv */
420     /* Flag PLAT_DR_FEATURE_ENABLED should only be set by DR driver. */
421     plat_dr_options &= ~PLAT_DR_FEATURE_ENABLED;
422 #ifndef __amd64
423     /* Only enable CPU/memory DR on 64 bits kernel. */
424     plat_dr_options &= ~PLAT_DR_FEATURE_MEMORY;
425     plat_dr_options &= ~PLAT_DR_FEATURE_CPU;
426 #endif /* __amd64 */
427 #endif /* __xpv */

429 /*
430  * Get value of "plat_dr_physmax" boot option.
431  * It overrides values calculated from MSCT or SRAT table.
432  */
433 if (bootprop_getval(PLAT_DR_PHYSMAX_NAME, &prop_value) == 0) {
434     plat_dr_physmax = ((uint64_t)prop_value) >> PAGESHIFT;
435 }

437 /* Get value of boot_ncpus. */
438 if (bootprop_getval(BOOT_NCPUS_NAME, &prop_value) != 0) {
439     boot_ncpus = NCPU;
440 } else {
441     boot_ncpus = (int)prop_value;
442     if (boot_ncpus <= 0 || boot_ncpus > NCPU)
443         boot_ncpus = NCPU;
444 }

446 /*
447  * Set max_ncpus and boot_max_ncpus to boot_ncpus if platform doesn't
448  * support CPU DR operations.
449  */
450 if (plat_dr_support_cpu() == 0) {
451     max_ncpus = boot_max_ncpus = boot_ncpus;
452 } else {
453     if (bootprop_getval(PLAT_MAX_NCPUS_NAME, &prop_value) != 0) {
454         max_ncpus = NCPU;
455     } else {
456         max_ncpus = (int)prop_value;
457         if (max_ncpus <= 0 || max_ncpus > NCPU) {
458             max_ncpus = NCPU;
459         }
460         if (boot_ncpus > max_ncpus) {
461             boot_ncpus = max_ncpus;
462         }
463     }
464 }

465 if (bootprop_getval(BOOT_MAX_NCPUS_NAME, &prop_value) != 0) {
466     boot_max_ncpus = boot_ncpus;
467 } else {
468     boot_max_ncpus = (int)prop_value;
469     if (boot_max_ncpus <= 0 || boot_max_ncpus > NCPU) {
470         boot_max_ncpus = boot_ncpus;
471     } else if (boot_max_ncpus > max_ncpus) {
472         boot_max_ncpus = max_ncpus;
473     }
474 }
475 }

477 /*
478  * Initialize the lgrp framework
479  */
480 lgrp_init(LGRP_INIT_STAGE1);

```

```
482     if (boothowto & RB_HALT) {
483         prom_printf("unix: kernel halted by -h flag\n");
484         prom_enter_mon();
485     }
487     ASSERT_STACK_ALIGNED();
489     /*
490     * Fill out cpu_ucode_info.  Update microcode if necessary.
491     */
492     ucode_check(CPU);
494     if (workaround_errata(CPU) != 0)
495         panic("critical workaround(s) missing for boot cpu");
496 }
```

unchanged portion omitted

\*\*\*\*\*

52277 Thu Jun 14 17:15:59 2018

new/usr/src/uts/i86pc/os/mp\_startup.c

9600 LDT still not happy under KPTI

\*\*\*\*\*

unchanged portion omitted

```
2047 /*
2048  * The following two routines are used as context operators on threads belonging
2049  * to processes with a private LDT (see sysi86). Due to the rarity of such
2050  * processes, these routines are currently written for best code readability and
2051  * organization rather than speed. We could avoid checking x86_featureset at
2052  * every context switch by installing different context ops, depending on
2053  * x86_featureset, at LDT creation time -- one for each combination of fast
2054  * syscall features.
2055  */
```

2057 /\*ARGSUSED\*/

2057 void

2058 **cpu\_fast\_syscall\_disable(void)**

2059 *cpu\_fast\_syscall\_disable(void \*arg)*

```
2059 {
2060     if (is_x86_feature(x86_featureset, X86FSET_MSR) &&
2061         is_x86_feature(x86_featureset, X86FSET_SEP))
2062         cpu_sep_disable();
2063     if (is_x86_feature(x86_featureset, X86FSET_MSR) &&
2064         is_x86_feature(x86_featureset, X86FSET_ASYSC))
2065         cpu_asysc_disable();
2066 }
```

2069 /\*ARGSUSED\*/

2068 void

2069 **cpu\_fast\_syscall\_enable(void)**

2071 *cpu\_fast\_syscall\_enable(void \*arg)*

```
2070 {
2071     if (is_x86_feature(x86_featureset, X86FSET_MSR) &&
2072         is_x86_feature(x86_featureset, X86FSET_SEP))
2073         cpu_sep_enable();
2074     if (is_x86_feature(x86_featureset, X86FSET_MSR) &&
2075         is_x86_feature(x86_featureset, X86FSET_ASYSC))
2076         cpu_asysc_enable();
2077 }
```

unchanged portion omitted



```

*****
38161 Thu Jun 14 17:15:59 2018
new/usr/src/uts/intel/ia32/os/desctbls.c
9600 LDT still not happy under KPTI
*****
_____unchanged_portion_omitted_____

166 /*
167  * The brand infrastructure interposes on two handlers, and we use one as a
168  * NULL signpost.
169  */
170 static struct interposing_handler brand_tbl[2];

172 /*
173  * software prototypes for default local descriptor table
174  */

176 /*
177  * Routines for loading segment descriptors in format the hardware
178  * can understand.
179  */

181 #if defined(__amd64)

181 /*
182  * In long mode we have the new L or long mode attribute bit
183  * for code segments. Only the conforming bit in type is used along
184  * with descriptor priority and present bits. Default operand size must
185  * be zero when in long mode. In 32-bit compatibility mode all fields
186  * are treated as in legacy mode. For data segments while in long mode
187  * only the present bit is loaded.
188  */
189 void
190 set_usegd(user_desc_t *dp, uint_t lmode, void *base, size_t size,
191          uint_t type, uint_t dpl, uint_t gran, uint_t defopsz)
192 {
193     ASSERT(lmode == SDP_SHORT || lmode == SDP_LONG);
194     /* This should never be a "system" segment. */
195     ASSERT3U(type & SDT_S, !=, 0);

197     /*
198     * 64-bit long mode.
199     */
200     if (lmode == SDP_LONG)
201         dp->usd_def32 = 0;          /* 32-bit operands only */
202     else
203         /*
204         * 32-bit compatibility mode.
205         */
206         dp->usd_def32 = defopsz;    /* 0 = 16, 1 = 32-bit ops */

208     /*
209     * We should always set the "accessed" bit (SDT_A), otherwise the CPU
210     * will write to the GDT whenever we change segment registers around.
211     * With KPTI on, the GDT is read-only in the user page table, which
212     * causes crashes if we don't set this.
213     */
214     ASSERT3U(type & SDT_A, !=, 0);

216     dp->usd_long = lmode;    /* 64-bit mode */
217     dp->usd_type = type;
218     dp->usd_dpl = dpl;
219     dp->usd_p = 1;
220     dp->usd_gran = gran;    /* 0 = bytes, 1 = pages */

222     dp->usd_lobase = (uintptr_t)base;

```

```

223     dp->usd_midbase = (uintptr_t)base >> 16;
224     dp->usd_hibase = (uintptr_t)base >> (16 + 8);
225     dp->usd_lolimit = size;
226     dp->usd_hilimit = (uintptr_t)size >> 16;
227 }

221 #elif defined(__i386)

229 /*
230  * Install user segment descriptor for code and data.
231  */
232 void
233 set_usegd(user_desc_t *dp, void *base, size_t size, uint_t type,
234          uint_t dpl, uint_t gran, uint_t defopsz)
235 {
236     dp->usd_lolimit = size;
237     dp->usd_hilimit = (uintptr_t)size >> 16;

239     dp->usd_lobase = (uintptr_t)base;
240     dp->usd_midbase = (uintptr_t)base >> 16;
241     dp->usd_hibase = (uintptr_t)base >> (16 + 8);
242 }

244 #endif /* __i386 */

246 /*
247  * Install system segment descriptor for LDT and TSS segments.
248  */

250 #if defined(__amd64)

253 void
254 set_syssegd(system_desc_t *dp, void *base, size_t size, uint_t type,
255            uint_t dpl)
256 {
257     dp->ssd_lolimit = size;
258     dp->ssd_hilimit = (uintptr_t)size >> 16;

260     dp->ssd_lobase = (uintptr_t)base;
261     dp->ssd_midbase = (uintptr_t)base >> 16;
262     dp->ssd_hibase = (uintptr_t)base >> (16 + 8);
263     dp->ssd_hi64base = (uintptr_t)base >> (16 + 8 + 8);

265     dp->ssd_type = type;
266     dp->ssd_zero1 = 0;        /* must be zero */
267     dp->ssd_zero2 = 0;
268     dp->ssd_dpl = dpl;
269     dp->ssd_p = 1;
270     dp->ssd_gran = 0;        /* force byte units */
271 }

_____unchanged_portion_omitted_____

284 #elif defined(__i386)

286 void
287 set_syssegd(system_desc_t *dp, void *base, size_t size, uint_t type,
288            uint_t dpl)
289 {
290     dp->ssd_lolimit = size;
291     dp->ssd_hilimit = (uintptr_t)size >> 16;

```

```

293     dp->ssd_lobase = (uintptr_t)base;
294     dp->ssd_midbase = (uintptr_t)base >> 16;
295     dp->ssd_hibase = (uintptr_t)base >> (16 + 8);

297     dp->ssd_type = type;
298     dp->ssd_zero = 0;      /* must be zero */
299     dp->ssd_dpl = dpl;
300     dp->ssd_p = 1;
301     dp->ssd_gran = 0;     /* force byte units */
302 }

304 void *
305 get_ssd_base(system_desc_t *dp)
306 {
307     uintptr_t     base;

309     base = (uintptr_t)dp->ssd_lobase |
310           (uintptr_t)dp->ssd_midbase << 16 |
311           (uintptr_t)dp->ssd_hibase << (16 + 8);
312     return ((void *)base);
313 }

315 #endif /* __i386 */

265 /*
266 * Install gate segment descriptor for interrupt, trap, call and task gates.
267 *
268 * For 64 bit native if we have KPTI enabled, we use the IST stack mechanism on
269 * all interrupts. We have different ISTs for each class of exceptions that are
270 * most likely to occur while handling an existing exception; while many of
271 * these are just going to panic, it's nice not to trample on the existing
272 * exception state for debugging purposes.
273 *
274 * Normal interrupts are all redirected unconditionally to the KPTI trampoline
275 * stack space. This unifies the trampoline handling between user and kernel
276 * space (and avoids the need to touch %gs).
277 *
278 * The KDI IDT *all* uses the DBG IST: consider single stepping tr_pftrap, when
279 * we do a read from KMDB that cause another #PF. Without its own IST, this
280 * would stomp on the kernel's mcpu_kpti_flt frame.
281 */
282 uint_t
283 idt_vector_to_ist(uint_t vector)
284 {
285 #if defined(__xpv)
286     _NOTE(ARGUNUSED(vector));
287     return (IST_NONE);
288 #else
289     switch (vector) {
290     /* These should always use IST even without KPTI enabled. */
291     case T_DBLFLT:
292         return (IST_DF);
293     case T_NMIFLT:
294         return (IST_NMI);
295     case T_MCE:
296         return (IST_MCE);

298     case T_BPTFLT:
299     case T_SGLSTP:
300         if (kpti_enable == 1) {
301             return (IST_DBG);
302         }
303         return (IST_NONE);
304     case T_STKFLT:
305     case T_GPFLT:

```

```

306     case T_PGFLT:
307         if (kpti_enable == 1) {
308             return (IST_NESTABLE);
309         }
310         return (IST_NONE);
311     default:
312         if (kpti_enable == 1) {
313             return (IST_DEFAULT);
314         }
315         return (IST_NONE);
316     }
317 #endif
318 }
    unchanged_portion_omitted

334 /*
335 * Updates a single user descriptor in the the GDT of the current cpu.
336 * Caller is responsible for preventing cpu migration.
337 */

339 void
340 gdt_update_usegd(uint_t sidx, user_desc_t *udp)
341 {
342 #if defined(DEBUG)
343     /* This should never be a "system" segment, but it might be null. */
344     if (udp->usd_p != 0 || udp->usd_type != 0) {
345         ASSERT3U(udp->usd_type & SDT_S, !=, 0);
346     }
347     /*
348      * We should always set the "accessed" bit (SDT_A), otherwise the CPU
349      * will write to the GDT whenever we change segment registers around.
350      * With KPTI on, the GDT is read-only in the user page table, which
351      * causes crashes if we don't set this.
352      */
353     if (udp->usd_p != 0 || udp->usd_type != 0) {
354         ASSERT3U(udp->usd_type & SDT_A, !=, 0);
355     }
356 #endif

358 #if defined(__xpv)
359     uint64_t dpa = CPU->cpu_m.mcpu_gdtpa + sizeof (*udp) * sidx;

361     if (HYPERVISOR_update_descriptor(pa_to_ma(dpa), *(uint64_t *)udp))
362         panic("gdt_update_usegd: HYPERVISOR_update_descriptor");

364 #else /* __xpv */
365     CPU->cpu_gdt[sidx] = *udp;
366 #endif /* __xpv */
367 }

369 /*
370 * Writes single descriptor pointed to by udp into a processes
371 * LDT entry pointed to by ldp.
372 */
373 int
374 ldt_update_segld(user_desc_t *ldp, user_desc_t *udp)
375 {
376 #if defined(DEBUG)
377     /* This should never be a "system" segment, but it might be null. */
378     if (udp->usd_p != 0 || udp->usd_type != 0) {
379         ASSERT3U(udp->usd_type & SDT_S, !=, 0);
380     }
381     /*

```

```
382  * We should always set the "accessed" bit (SDT_A), otherwise the CPU
383  * will write to the LDT whenever we change segment registers around.
384  * With KPTI on, the LDT is read-only in the user page table, which
385  * causes crashes if we don't set this.
386  */
387  if (udp->usd_p != 0 || udp->usd_type != 0) {
388      ASSERT3U(udp->usd_type & SDT_A, !=, 0);
389  }
390  #endif

392 #if defined(__xpv)

393     uint64_t dpa;

395     dpa = mmu_ptob(hat_getpfnum(kas.a_hat, (caddr_t)ldp) |
396                 ((uintptr_t)ldp & PAGEOFFSET);

398     /*
399     * The hypervisor is a little more restrictive about what it
400     * supports in the LDT.
401     */
402     if (HYPERVISOR_update_descriptor(pa_to_ma(dpa), *(uint64_t *)udp) != 0)
403         return (EINVAL);

405 #else /* __xpv */

406     *ldp = *udp;

408 #endif /* __xpv */
409     return (0);
410 }
unchanged_portion_omitted
```

```
*****
```

```
21632 Thu Jun 14 17:15:59 2018
```

```
new/usr/src/uts/intel/ia32/os/sysi86.c
```

```
9600 LDT still not happy under KPTI
```

```
*****
```

```
_____ unchanged_portion_omitted _____
```

```
277 static void
278 ssd_to_usd(struct ssd *ssd, user_desc_t *usd)
279 {
281     ASSERT(bcmp(usd, &null_udesc, sizeof (*usd)) == 0);
283     USEGD_SETBASE(usd, ssd->bo);
284     USEGD_SETLIMIT(usd, ssd->ls);
286     /*
287      * Set type, dpl and present bits.
288      *
289      * Force the "accessed" bit to on so that we don't run afoul of
290      * KPTI.
291      * set type, dpl and present bits.
292      */
292     usd->usd_type = ssd->acc1 | SDT_A;
289     usd->usd_type = ssd->acc1;
293     usd->usd_dpl = ssd->acc1 >> 5;
294     usd->usd_p = ssd->acc1 >> (5 + 2);
296     ASSERT(usd->usd_type >= SDT_MEMRO);
297     ASSERT(usd->usd_dpl == SEL_UPL);
299     /*
300      * 64-bit code selectors are never allowed in the LDT.
301      * Reserved bit is always 0 on 32-bit systems.
302      */
303 #if defined(__amd64)
304     usd->usd_long = 0;
305 #else
306     usd->usd_reserved = 0;
307 #endif
309     /*
310      * set avl, DB and granularity bits.
311      */
312     usd->usd_avl = ssd->acc2;
313     usd->usd_def32 = ssd->acc2 >> (1 + 1);
314     usd->usd_gran = ssd->acc2 >> (1 + 1 + 1);
315 }
_____ unchanged_portion_omitted _____
342 #endif /* __i386 */
344 /*
345  * Load LDT register with the current process's LDT.
346  */
347 static void
348 ldt_load(void)
349 {
350 #if defined(__xpv)
351     xen_set_ldt(curproc->p_ldt, curproc->p_ldtlimit + 1);
348     xen_set_ldt(get_ssd_base(&curproc->p_ldt_desc),
349     curproc->p_ldtlimit + 1);
352 #else
353     size_t len;
354     system_desc_t desc;
```

```
356     /*
357      * Before we can use the LDT on this CPU, we must install the LDT in the
358      * user mapping table.
359      */
360     len = (curproc->p_ldtlimit + 1) * sizeof (user_desc_t);
361     bcopy(curproc->p_ldt, CPU->cpu_m.mcpu_ldt, len);
362     CPU->cpu_m.mcpu_ldt_len = len;
363     set_syssegd(&desc, CPU->cpu_m.mcpu_ldt, len - 1, SDT_SYSLDT, SEL_KPL);
364     *((system_desc_t *)&CPU->cpu_gdt[GDT_LDT]) = desc;
366     wr_ldtr(ULDT_SEL);
367 #endif
368 }
_____ unchanged_portion_omitted _____
388 /*ARGSUSED*/
389 static void
390 ldt_savectx(proc_t *p)
391 {
392     ASSERT(p->p_ldt != NULL);
393     ASSERT(p == curproc);
395 #if defined(__amd64)
396     /*
397      * The 64-bit kernel must be sure to clear any stale ldt
398      * selectors when context switching away from a process that
399      * has a private ldt. Consider the following example:
400      *
401      * Wine creates a ldt descriptor and points a segment register
402      * to it.
403      *
404      * We then context switch away from wine lwp to kernel
405      * thread and hit breakpoint in kernel with kmdb
406      *
407      * When we continue and resume from kmdb we will #gp
408      * fault since kmdb will have saved the stale ldt selector
409      * from wine and will try to restore it but we are no longer in
410      * the context of the wine process and do not have our
411      * ldtr register pointing to the private ldt.
412      */
413     reset_sregs();
414 #endif
416     ldt_unload();
417     cpu_fast_syscall_enable();
415     cpu_fast_syscall_enable(NULL);
418 }
420 static void
421 ldt_restorectx(proc_t *p)
422 {
423     ASSERT(p->p_ldt != NULL);
424     ASSERT(p == curproc);
426     ldt_load();
427     cpu_fast_syscall_disable();
425     cpu_fast_syscall_disable(NULL);
428 }
430 /*
431  * At exec time, we need to clear up our LDT context and re-enable fast syscalls
432  * for the new process image.
433  *
434  * The same is true for the other case, where we have:
435  *
436  * proc_exit()
```

```

437 * ->exitpctx()->ldt_savectx()
438 * ->freepctx()->ldt_freectx()
439 *
440 * Because pre-emption is not prevented between the two callbacks, we could have
441 * come off CPU, and brought back LDT context when coming back on CPU via
442 * ldt_restorectx().
443 * When a process with a private LDT execs, fast syscalls must be enabled for
444 * the new process image.
445 */
446 /* ARGSUSED */
447 static void
448 ldt_freectx(proc_t *p, int isexec)
449 {
450     ASSERT(p->p_ldt != NULL);
451     ASSERT(p == curproc);
452     ASSERT(p->p_ldt);
453
454     if (isexec) {
455         kpreempt_disable();
456         ldt_free(p);
457         cpu_fast_syscall_enable();
458         cpu_fast_syscall_enable(NULL);
459         kpreempt_enable();
460     }
461
462     /*
463      * ldt_free() will free the memory used by the private LDT, reset the
464      * process's descriptor, and re-program the LDTR.
465      */
466     ldt_free(p);
467 }
468
469 #ifndef unchanged_portion_omitted_
470
471 int
472 setdscr(struct ssd *ssd)
473 {
474     ushort_t seli;          /* selector index */
475     user_desc_t *ldp;       /* descriptor pointer */
476     user_desc_t ndesc;     /* new descriptor */
477     proc_t *pp = curproc;
478     proc_t *pp = ttoproc(curthread);
479     int rc = 0;
480
481     /*
482      * LDT segments: executable and data at DPL 3 only.
483      */
484     if (!SELISLDT(ssd->sel) || !SELISUPL(ssd->sel))
485         return (EINVAL);
486
487     /*
488      * check the selector index.
489      */
490     seli = SELTOIDX(ssd->sel);
491     if (seli >= MAXNLDT || seli < LDT_UBBASE)
492         return (EINVAL);
493
494     ndesc = null_udesc;
495     mutex_enter(&pp->p_ldtlock);
496
497     /*
498      * If this is the first time for this process then setup a
499      * private LDT for it.
500      */
501     if (pp->p_ldt == NULL) {
502         ldt_alloc(pp, seli);
503     }
504 }
505

```

```

535     /*
536      * Now that this process has a private LDT, the use of
537      * the syscall/sysret and sysenter/sysexit instructions
538      * is forbidden for this processes because they destroy
539      * the contents of %cs and %ss segment registers.
540      *
541      * Explicitly disable them here and add a context handler
542      * to the process. Note that disabling
543      * them here means we can't use sysret or sysexit on
544      * the way out of this system call - so we force this
545      * thread to take the slow path (which doesn't make use
546      * of sysenter or sysexit) back out.
547      */
548     kpreempt_disable();
549     ldt_installctx(pp, NULL);
550     cpu_fast_syscall_disable();
551     cpu_fast_syscall_disable(NULL);
552     ASSERT(curthread->t_post_sys != 0);
553     kpreempt_enable();
554
555     } else if (seli > pp->p_ldtlimit) {
556         ASSERT(pp->p_pctx != NULL);
557
558         /*
559          * Increase size of ldt to include seli.
560          */
561         ldt_grow(pp, seli);
562     }
563
564     ASSERT(seli <= pp->p_ldtlimit);
565     ldp = &pp->p_ldt[seli];
566
567     /*
568      * On the 64-bit kernel, this is where things get more subtle.
569      * Recall that in the 64-bit kernel, when we enter the kernel we
570      * deliberately -don't- reload the segment selectors we came in on
571      * for %ds, %es, %fs or %gs. Messing with selectors is expensive,
572      * and the underlying descriptors are essentially ignored by the
573      * hardware in long mode - except for the base that we override with
574      * the gsbase MSRs.
575      *
576      * However, there's one unfortunate issue with this rosy picture --
577      * a descriptor that's not marked as 'present' will still generate
578      * an #np when loading a segment register.
579      *
580      * Consider this case. An lwp creates a harmless LDT entry, points
581      * one of it's segment registers at it, then tells the kernel (here)
582      * to delete it. In the 32-bit kernel, the #np will happen on the
583      * way back to userland where we reload the segment registers, and be
584      * handled in kern_gpfault(). In the 64-bit kernel, the same thing
585      * will happen in the normal case too. However, if we're trying to
586      * use a debugger that wants to save and restore the segment registers,
587      * and the debugger things that we have valid segment registers, we
588      * have the problem that the debugger will try and restore the
589      * segment register that points at the now 'not present' descriptor
590      * and will take a #np right there.
591      *
592      * We should obviously fix the debugger to be paranoid about
593      * -not- restoring segment registers that point to bad descriptors;
594      * however we can prevent the problem here if we check to see if any
595      * of the segment registers are still pointing at the thing we're
596      * destroying; if they are, return an error instead. (That also seems
597      * a lot better failure mode than SIGKILL and a core file
598      * from kern_gpfault() too.)
599      */
600     if (SI86SSD_PRESENT(ssd) == 0) {

```

```

600     kthread_t *t;
601     int bad = 0;

603     /*
604     * Look carefully at the segment registers of every lwp
605     * in the process (they're all stopped by our caller).
606     * If we're about to invalidate a descriptor that's still
607     * being referenced by *any* of them, return an error,
608     * rather than having them #gp on their way out of the kernel.
609     */
610     ASSERT(pp->p_lwprcnt == 1);

612     mutex_enter(&pp->p_lock);
613     t = pp->p_tlist;
614     do {
615         klwp_t *lwp = ttolwp(t);
616         struct regs *rp = lwp->lwp_regs;
617 #if defined(__amd64)
618         pcb_t *pcb = &lwp->lwp_pcb;
619 #endif

621         if (ssd->sel == rp->r_cs || ssd->sel == rp->r_ss) {
622             bad = 1;
623             break;
624         }

626 #if defined(__amd64)
627         if (pcb->pcb_rupdate == 1) {
628             if (ssd->sel == pcb->pcb_ds ||
629                 ssd->sel == pcb->pcb_es ||
630                 ssd->sel == pcb->pcb_fs ||
631                 ssd->sel == pcb->pcb_gs) {
632                 bad = 1;
633                 break;
634             }
635         } else
636 #endif
637         {
638             if (ssd->sel == rp->r_ds ||
639                 ssd->sel == rp->r_es ||
640                 ssd->sel == rp->r_fs ||
641                 ssd->sel == rp->r_gs) {
642                 bad = 1;
643                 break;
644             }
645         }

647     } while ((t = t->t_forw) != pp->p_tlist);
648     mutex_exit(&pp->p_lock);

650     if (bad) {
651         mutex_exit(&pp->p_ldtlock);
652         return (EBUSY);
653     }
654 }

656 /*
657 * If accl is zero, clear the descriptor (including the 'present' bit).
658 * Make sure we update the CPU-private copy of the LDT.
659 * If accl is zero, clear the descriptor (including the 'present' bit)
660 */
661 if (ssd->accl == 0) {
662     rc = ldt_update_segdesc(ldp, &null_udesc);
663     kpreempt_disable();
664     ldt_load();
665     kpreempt_enable();

```

```

665         mutex_exit(&pp->p_ldtlock);
666         return (rc);
667     }

669     /*
670     * Check segment type, allow segment not present and
671     * only user DPL (3).
672     */
673     if (SI86SSD_DPL(ssd) != SEL_UPL) {
674         mutex_exit(&pp->p_ldtlock);
675         return (EINVAL);
676     }

667 #if defined(__amd64)
678     /*
679     * Do not allow 32-bit applications to create 64-bit mode code
680     * segments.
681     */
682     if (SI86SSD_ISUSEG(ssd) && ((SI86SSD_TYPE(ssd) >> 3) & 1) == 1 &&
683         SI86SSD_ISLONG(ssd)) {
684         mutex_exit(&pp->p_ldtlock);
685         return (EINVAL);
686     }
687 #endif /* __amd64 */

688     /*
689     * Set up a code or data user segment descriptor, making sure to update
690     * the CPU-private copy of the LDT.
691     * Set up a code or data user segment descriptor.
692     */
693     if (SI86SSD_ISUSEG(ssd)) {
694         ssd_to_usd(ssd, &ndesc);
695         rc = ldt_update_segdesc(ldp, &ndesc);
696         kpreempt_disable();
697         ldt_load();
698         kpreempt_enable();
699         mutex_exit(&pp->p_ldtlock);
700         return (rc);
701     }

689 #if defined(__i386)
690     /*
691     * Allow a call gate only if the destination is in the LDT
692     * and the system is running in 32-bit legacy mode.
693     *
694     * In long mode 32-bit call gates are redefined as 64-bit call
695     * gates and the hw enforces that the target code selector
696     * of the call gate must be 64-bit selector. A #gp fault is
697     * generated if otherwise. Since we do not allow 32-bit processes
698     * to switch themselves to 64-bits we never allow call gates
699     * on 64-bit system system.
700     */
701     if (SI86SSD_TYPE(ssd) == SGT_SYSCGT && SELISLDT(ssd->ls)) {

704         ssd_to_sgd(ssd, (gate_desc_t *)&ndesc);
705         rc = ldt_update_segdesc(ldp, &ndesc);
706         mutex_exit(&pp->p_ldtlock);
707         return (rc);
708     }
709 #endif /* __i386 */

711     mutex_exit(&pp->p_ldtlock);
712     return (EINVAL);
713 }

```

```

706 /*
707 * Allocate new LDT for process just large enough to contain seli. Note we
708 * allocate and grow LDT in PAGESIZE chunks. We do this to simplify the
709 * implementation and because on the hypervisor it's required, since the LDT
710 * must live on pages that have PROT_WRITE removed and which are given to the
711 * hypervisor.
712 *
713 * Note that we don't actually load the LDT into the current CPU here: it's done
714 * later by our caller.
715 */
716 static void
717 ldt_alloc(proc_t *pp, uint_t seli)
718 {
719     user_desc_t    *ldt;
720     size_t         ldtasz;
721     uint_t         nsels;
722
723     ASSERT(MUTEX_HELD(&pp->p_ldtlock));
724     ASSERT(pp->p_ldt == NULL);
725     ASSERT(pp->p_ldtlimit == 0);
726
727     /*
728     * Allocate new LDT just large enough to contain seli. The LDT must
729     * always be allocated in units of pages for KPTI.
730     */
731     ldtasz = P2ROUNDUP((seli + 1) * sizeof (user_desc_t), PAGESIZE);
732     nsels = ldtasz / sizeof (user_desc_t);
733     ASSERT(nsels >= MINNLDT && nsels <= MAXNLDT);
734
735     ldt = kmem_zalloc(ldtasz, KM_SLEEP);
736     ASSERT(IS_P2ALIGNED(ldt, PAGESIZE));
737
738 #if defined(__xpv)
739     if (xen_ldt_setprot(ldt, ldtasz, PROT_READ))
740         panic("ldt_alloc:xen_ldt_setprot(PROT_READ) failed");
741 #endif
742
743     pp->p_ldt = ldt;
744     pp->p_ldtlimit = nsels - 1;
745     set_syssegd(&pp->p_ldt_desc, ldt, ldtasz - 1, SDT_SYSLDT, SEL_KPL);
746
747     if (pp == curproc) {
748         kpreempt_disable();
749         ldt_load();
750         kpreempt_enable();
751     }
752 }
753
754 static void
755 ldt_free(proc_t *pp)
756 {
757     user_desc_t    *ldt;
758     size_t         ldtasz;
759
760     ASSERT(pp->p_ldt != NULL);
761
762     mutex_enter(&pp->p_ldtlock);
763     ldt = pp->p_ldt;
764     ldtasz = (pp->p_ldtlimit + 1) * sizeof (user_desc_t);
765
766     ASSERT(IS_P2ALIGNED(ldt, PAGESIZE));

```

```

761     pp->p_ldt = NULL;
762     pp->p_ldtlimit = 0;
763     pp->p_ldt_desc = null_sdesc;
764     mutex_exit(&pp->p_ldtlock);
765
766     if (pp == curproc) {
767         kpreempt_disable();
768         ldt_unload();
769         kpreempt_enable();
770     }
771 #if defined(__xpv)
772     /*
773     * We are not allowed to make the ldt writable until after
774     * we tell the hypervisor to unload it.
775     */
776     if (xen_ldt_setprot(ldt, ldtasz, PROT_READ | PROT_WRITE))
777         panic("ldt_free:xen_ldt_setprot(PROT_READ|PROT_WRITE) failed");
778 #endif
779
780     kmem_free(ldt, ldtasz);
781 }
782
783 unchanged portion omitted
784
785 /*
786 * Note that we don't actually load the LDT into the current CPU here: it's done
787 * later by our caller - unless we take an error. This works out because
788 * ldt_load() does a copy of ->p_ldt instead of directly loading it into the GDT
789 * (and therefore can't be using the freed old LDT), and by definition if the
790 * new entry didn't pass validation, then the proc shouldn't be referencing an
791 * entry in the extended region.
792 */
793 static void
794 ldt_grow(proc_t *pp, uint_t seli)
795 {
796     user_desc_t    *oldt, *nldt;
797     uint_t         nsels;
798     size_t         oldtsz, nldtsz;
799
800     ASSERT(MUTEX_HELD(&pp->p_ldtlock));
801     ASSERT(pp->p_ldt != NULL);
802     ASSERT(pp->p_ldtlimit != 0);
803
804     /*
805     * Allocate larger LDT just large enough to contain seli. The LDT must
806     * always be allocated in units of pages for KPTI.
807     */
808     nldtsz = P2ROUNDUP((seli + 1) * sizeof (user_desc_t), PAGESIZE);
809     nsels = nldtsz / sizeof (user_desc_t);
810     ASSERT(nsels >= MINNLDT && nsels <= MAXNLDT);
811     ASSERT(nsels > pp->p_ldtlimit);
812
813     oldt = pp->p_ldt;
814     oldtsz = (pp->p_ldtlimit + 1) * sizeof (user_desc_t);
815
816     nldt = kmem_zalloc(nldtsz, KM_SLEEP);
817     ASSERT(IS_P2ALIGNED(nldt, PAGESIZE));
818
819     bcopy(oldt, nldt, oldtsz);
820
821     /*
822     * unload old ldt.
823     */
824     kpreempt_disable();
825     ldt_unload();

```

```
868     kpreempt_enable();

870 #if defined(__xpv)

872     /*
873     * Make old ldt writable and new ldt read only.
874     */
875     if (xen_ldt_setprot(oldt, oldtsz, PROT_READ | PROT_WRITE))
876         panic("ldt_grow:xen_ldt_setprot(PROT_READ|PROT_WRITE) failed");

878     if (xen_ldt_setprot(nldt, nldtsz, PROT_READ))
879         panic("ldt_grow:xen_ldt_setprot(PROT_READ) failed");
880 #endif

882     pp->p_ldt = nldt;
883     pp->p_ldtlimit = nsels - 1;

891     /*
892     * write new ldt segment descriptor.
893     */
894     set_syssegd(&pp->p_ldt_desc, nldt, nldtsz - 1, SDT_SYSLDT, SEL_KPL);

896     /*
897     * load the new ldt.
898     */
899     kpreempt_disable();
900     ldt_load();
901     kpreempt_enable();

885     kmem_free(oldt, oldtsz);
886 }
unchanged_portion_omitted
```



```

*****
25075 Thu Jun 14 17:16:00 2018
new/usr/src/uts/intel/sys/segments.h
9600 LDT still not happy under KPFI
*****
_____unchanged_portion_omitted_____

384 #define GATESEG_GETOFFSET(sgd) ((uintptr_t)((sgd)->sgd_loffset | \
385 (sgd)->sgd_hioffset << 16 | \
386 (uint64_t)((sgd)->sgd_hi64offset) << 32))

388 #endif /* __amd64 */

390 /*
391 * functions for initializing and updating segment descriptors.
392 */
393 #if defined(__amd64)

395 extern void set_usegd(user_desc_t *, uint_t, void *, size_t, uint_t, uint_t,
396 uint_t, uint_t);

398 #elif defined(__i386)

400 extern void set_usegd(user_desc_t *, void *, size_t, uint_t, uint_t,
401 uint_t, uint_t);

403 #endif /* __i386 */

405 extern uint_t idt_vector_to_ist(uint_t);

407 extern void set_gatesegd(gate_desc_t *, void (*)(void), selector_t,
408 uint_t, uint_t, uint_t);

410 extern void set_sysseg(system_desc_t *, void *, size_t, uint_t, uint_t);

412 extern void *get_ssd_base(system_desc_t *);

414 extern void gdt_update_usegd(uint_t, user_desc_t *);

416 extern int ldt_update_seg(system_desc_t *, user_desc_t *);

418 #if defined(__xpv)

420 extern int xen_idt_to_trap_info(uint_t, gate_desc_t *, void *);
421 extern void xen_idt_write(gate_desc_t *, uint_t);

423 #endif /* __xen */

425 void init_boot_gdt(user_desc_t *);

427 #endif /* _ASM */

429 /*
430 * Common segment parameter definitions for granularity, default
431 * operand size and operation mode.
432 */
433 #define SDP_BYTES 0 /* segment limit scaled to bytes */
434 #define SDP_PAGES 1 /* segment limit scaled to pages */
435 #define SDP_OP32 1 /* code and data default operand = 32 bits */
436 #define SDP_LONG 1 /* long mode code segment (64 bits) */
437 #define SDP_SHORT 0 /* compat/legacy code segment (32 bits) */
438 /*
439 * System segments and gate types.
440 *
441 * In long mode i386 32-bit ldt, tss, call, interrupt and trap gate
442 * types are redefined into 64-bit equivalents.

```

```

443 */
444 #define SDT_SYSNULL 0 /* system null */
445 #define SDT_SYS286TSS 1 /* system 286 TSS available */
446 #define SDT_SYS286TSS 2 /* system local descriptor table */
447 #define SDT_SYS286BSY 3 /* system 286 TSS busy */
448 #define SDT_SYS286CGT 4 /* system 286 call gate */
449 #define SDT_SYSTASKGT 5 /* system task gate */
450 #define SDT_SYS286IGT 6 /* system 286 interrupt gate */
451 #define SDT_SYS286TGT 7 /* system 286 trap gate */
452 #define SDT_SYSNULL2 8 /* system null again */
453 #define SDT_SYSTSS 9 /* system TSS available */
454 #define SDT_SYSNULL3 10 /* system null again */
455 #define SDT_SYSTSSBSY 11 /* system TSS busy */
456 #define SDT_SYSCGT 12 /* system call gate */
457 #define SDT_SYSNULL4 13 /* system null again */
458 #define SDT_SYSIGT 14 /* system interrupt gate */
459 #define SDT_SYSTGT 15 /* system trap gate */

461 /*
462 * Memory segment types.
463 *
464 * While in long mode expand-down, writable and accessed type field
465 * attributes are ignored. Only the conforming bit is loaded by hardware
466 * for long mode code segment descriptors.
467 */
468 #define SDT_MEMRO 16 /* read only */
469 #define SDT_MEMROA 17 /* read only accessed */
470 #define SDT_MEMRW 18 /* read write */
471 #define SDT_MEMRWA 19 /* read write accessed */
472 #define SDT_MEMROD 20 /* read only expand dwn limit */
473 #define SDT_MEMRODA 21 /* read only expand dwn limit accessed */
474 #define SDT_MEMRWD 22 /* read write expand dwn limit */
475 #define SDT_MEMRWDA 23 /* read write expand dwn limit accessed */
476 #define SDT_MEME 24 /* execute only */
477 #define SDT_MEMEA 25 /* execute only accessed */
478 #define SDT_MEMER 26 /* execute read */
479 #define SDT_MEMERA 27 /* execute read accessed */
480 #define SDT_MEMEAC 28 /* execute only conforming */
481 #define SDT_MEMEAC 29 /* execute only accessed conforming */
482 #define SDT_MEMERC 30 /* execute read conforming */
483 #define SDT_MEMERAC 31 /* execute read accessed conforming */

485 /* These bits are within the "type" field, like the values above. */
486 #define SDT_A 0x01 /* accessed bit */
487 #define SDT_S 0x10 /* S-bit at the top of "type" for usegs */

489 /*
490 * Entries in the Interrupt Descriptor Table (IDT)
491 */
492 #define IDT_DE 0 /* #DE: Divide Error */
493 #define IDT_DB 1 /* #DB: Debug */
494 #define IDT_NMI 2 /* Nonmaskable External Interrupt */
495 #define IDT_BP 3 /* #BP: Breakpoint */
496 #define IDT_OF 4 /* #OF: Overflow */
497 #define IDT_BR 5 /* #BR: Bound Range Exceeded */
498 #define IDT_UD 6 /* #UD: Undefined/Invalid Opcode */
499 #define IDT_NM 7 /* #NM: No Math Coprocessor */
500 #define IDT_DF 8 /* #DF: Double Fault */
501 #define IDT_FPUGP 9 /* Coprocessor Segment Overrun */
502 #define IDT_TS 10 /* #TS: Invalid TSS */
503 #define IDT_NP 11 /* #NP: Segment Not Present */
504 #define IDT_SS 12 /* #SS: Stack Segment Fault */
505 #define IDT_GP 13 /* #GP: General Protection Fault */
506 #define IDT_PF 14 /* #PF: Page Fault */
507 #define IDT_MF 16 /* #MF: FPU Floating-Point Error */
508 #define IDT_AC 17 /* #AC: Alignment Check */

```

```

509 #define IDT_MC          18      /* #MC: Machine Check */
510 #define IDT_XF          19      /* #XF: SIMD Floating-Point Exception */
511 #define NIDT            256     /* size in entries of IDT */

513 /*
514 * Entries in the Global Descriptor Table (GDT)
515 *
516 * We make sure to space the system descriptors (LDT's, TSS')
517 * such that they are double gdt slot aligned. This is because
518 * in long mode system segment descriptors expand to 128 bits.
519 *
520 * GDT_LWPFs and GDT_LWPGs must be the same for both 32 and 64-bit
521 * kernels. See setup_context in libc. 64-bit processes must set
522 * %fs or %gs to null selector to use 64-bit fsbase or gsbase
523 * respectively.
524 */
525 #define GDT_NULL        0        /* null */
526 #define GDT_B32DATA    1        /* dboot 32 bit data descriptor */
527 #define GDT_B32CODE    2        /* dboot 32 bit code descriptor */
528 #define GDT_B16CODE    3        /* bios call 16 bit code descriptor */
529 #define GDT_B16DATA    4        /* bios call 16 bit data descriptor */
530 #define GDT_B64CODE    5        /* dboot 64 bit code descriptor */
531 #define GDT_BGSTMP     7        /* kmdb descriptor only used early in boot */
532 #define GDT_CPUID      16       /* store numeric id of current CPU */

534 #if defined(__amd64)

536 #define GDT_KCODE      6        /* kernel code seg %cs */
537 #define GDT_KDATA      7        /* kernel data seg %ds */
538 #define GDT_U32CODE    8        /* 32-bit process on 64-bit kernel %cs */
539 #define GDT_UDATA      9        /* user data seg %ds (32 and 64 bit) */
540 #define GDT_UCODE      10       /* native user code seg %cs */
541 #define GDT_LDT        12       /* (12-13) LDT for current process */
542 #define GDT_KTSS       14       /* (14-15) kernel tss */
543 #define GDT_FS         GDT_NULL /* kernel %fs segment selector */
544 #define GDT_GS         GDT_NULL /* kernel %gs segment selector */
545 #define GDT_LWPFs     55       /* lwp private %fs segment selector (32-bit) */
546 #define GDT_LWPGs     56       /* lwp private %gs segment selector (32-bit) */
547 #define GDT_BRANDMIN   57       /* first entry in GDT for brand usage */
548 #define GDT_BRANDMAX   61       /* last entry in GDT for brand usage */
549 #define NGDT           62       /* number of entries in GDT */

551 /*
552 * This selector is only used in the temporary GDT used to bring additional
553 * CPUs from 16-bit real mode into long mode in real_mode_start().
554 */
555 #define TEMP_GDT_KCODE64 1       /* 64-bit code selector */

557 #elif defined(__i386)

559 #define GDT_LDT        40       /* LDT for current process */
560 #define GDT_KTSS       42       /* kernel tss */
561 #define GDT_KCODE      43       /* kernel code seg %cs */
562 #define GDT_KDATA      44       /* kernel data seg %ds */
563 #define GDT_UCODE      45       /* native user code seg %cs */
564 #define GDT_UDATA      46       /* user data seg %ds (32 and 64 bit) */
565 #define GDT_DBFLT      47       /* double fault #DF selector */
566 #define GDT_FS         53       /* kernel %fs segment selector */
567 #define GDT_GS         54       /* kernel %gs segment selector */
568 #define GDT_LWPFs     55       /* lwp private %fs segment selector */
569 #define GDT_LWPGs     56       /* lwp private %gs segment selector */
570 #define GDT_BRANDMIN   57       /* first entry in GDT for brand usage */
571 #define GDT_BRANDMAX   61       /* last entry in GDT for brand usage */
572 #if !defined(__xpv)
573 #define NGDT           90       /* number of entries in GDT */
574 #else

```

```

575 #define NGDT           512     /* single 4K page for the hypervisor */
576 #endif

578 #endif /* __i386 */

580 /*
581 * Convenient selector definitions.
582 */

584 /*
585 * XXPV 64 bit Xen only allows the guest %cs/%ss be the private ones it
586 * provides, not the ones we create for ourselves. See FLAT_RING3_CS64 in
587 * public/arch-x86_64.h
588 *
589 * 64-bit Xen runs paravirtual guests in ring 3 but emulates them running in
590 * ring 0 by clearing CPL in %cs value pushed on guest exception stacks.
591 * Therefore we will have KCS_SEL value indicate ring 0 and use that everywhere
592 * in the kernel. But in the few files where we initialize segment registers or
593 * create and update descriptors we will explicitly OR in SEL_KPL (ring 3) for
594 * kernel %cs. See desctbls.c for an example.
595 */

597 #if defined(__xpv) && defined(__amd64)
598 #define KCS_SEL        0xe030    /* FLAT_RING3_CS64 & 0xFFFF0 */
599 #define KDS_SEL        0xe02b    /* FLAT_RING3_SS64 */
600 #else
601 #define KCS_SEL        SEL_GDT(GDT_KCODE, SEL_KPL)
602 #define KDS_SEL        SEL_GDT(GDT_KDATA, SEL_KPL)
603 #endif

605 #define UCS_SEL        SEL_GDT(GDT_UCODE, SEL_UPL)
606 #if defined(__amd64)
607 #define TEMP_CS64_SEL  SEL_GDT(TEMP_GDT_KCODE64, SEL_KPL)
608 #define U32CS_SEL     SEL_GDT(GDT_U32CODE, SEL_UPL)
609 #endif

611 #define UDS_SEL        SEL_GDT(GDT_UDATA, SEL_UPL)
612 #define ULDT_SEL       SEL_GDT(GDT_LDT, SEL_KPL)
613 #define KTSS_SEL       SEL_GDT(GDT_KTSS, SEL_KPL)
614 #define DFTSS_SEL      SEL_GDT(GDT_DBFLT, SEL_KPL)
615 #define KFS_SEL        0
616 #define KGS_SEL        SEL_GDT(GDT_GS, SEL_KPL)
617 #define LWPFs_SEL      SEL_GDT(GDT_LWPFs, SEL_UPL)
618 #define LWPGs_SEL      SEL_GDT(GDT_LWPGs, SEL_UPL)
619 #define BRANDMIN_SEL   SEL_GDT(GDT_BRANDMIN, SEL_UPL)
620 #define BRANDMAX_SEL   SEL_GDT(GDT_BRANDMAX, SEL_UPL)

622 #define B64CODE_SEL    SEL_GDT(GDT_B64CODE, SEL_KPL)
623 #define B32CODE_SEL    SEL_GDT(GDT_B32CODE, SEL_KPL)
624 #define B32DATA_SEL    SEL_GDT(GDT_B32DATA, SEL_KPL)
625 #define B16CODE_SEL    SEL_GDT(GDT_B16CODE, SEL_KPL)
626 #define B16DATA_SEL    SEL_GDT(GDT_B16DATA, SEL_KPL)

628 /*
629 * Temporary %gs descriptor used by kmdb with -d option. Only lives
630 * in boot's GDT and is not copied into kernel's GDT from boot.
631 */
632 #define KMDBGs_SEL     SEL_GDT(GDT_BGSTMP, SEL_KPL)

634 /*
635 * Selector used for kdi_idt when kmdb has taken over the IDT.
636 */
637 #if defined(__amd64)
638 #define KMDBCODE_SEL    B64CODE_SEL
639 #else
640 #define KMDBCODE_SEL    B32CODE_SEL

```

```

641 #endif

643 /*
644  * Entries in default Local Descriptor Table (LDT) for every process.
645  */
646 #define LDT_SYSCALL      0      /* call gate for libc.a (obsolete) */
647 #define LDT_SIGCALL     1      /* EOL me, call gate for static sigreturn */
648 #define LDT_RESVD1      2      /* old user %cs */
649 #define LDT_RESVD2      3      /* old user %ds */
650 #define LDT_ALTSYSCALL  4      /* alternate call gate for system calls */
651 #define LDT_ALTSIGCALL  5      /* EOL me, alternate call gate for sigreturn */
652 #define LDT_UDBASE      6      /* user descriptor base index */
653 #define MINNLDT        512     /* Current min solaris ldt size (1 4K page) */
654 #define MAXNLDT        8192    /* max solaris ldt size (16 4K pages) */

656 #ifdef _KERNEL
657 #define LDT_CPU_SIZE    (16 * 4096) /* Size of kernel per-CPU allocation */
658 #endif

660 #ifndef _ASM

662 extern gate_desc_t      *idt0;
663 extern desc_tbr_t      idt0_default_reg;
664 extern user_desc_t     *gdt0;

666 extern user_desc_t     zero_udesc;
667 extern user_desc_t     null_udesc;
668 extern system_desc_t   null_sdesc;

670 #if defined(__amd64)
671 extern user_desc_t     zero_u32desc;
672 #endif
673 #if defined(__amd64)
674 extern user_desc_t     ucs_on;
675 extern user_desc_t     ucs_off;
676 extern user_desc_t     ucs32_on;
677 extern user_desc_t     ucs32_off;
678 #endif /* __amd64 */

680 extern tss_t *ktss0;

682 #if defined(__i386)
683 extern tss_t *dftss0;
684 #endif /* __i386 */

686 extern void div0trap(), dbgtrap(), nmiint(), brktrap(), ovflotrap();
687 extern void boundstrap(), invoptrap(), ndptrap();
688 #if !defined(__xpv)
689 extern void syserrtrap();
690 #endif
691 extern void invaltrap(), invtsstrap(), segnptrap(), stktrap();
692 extern void gptrap(), pfttrap(), ndperr();
693 extern void overrun(), resvtrap();
694 extern void _start(), cmnint();
695 extern void achktrap(), mcetrap();
696 extern void xmtrap();
697 extern void fasttrap();
698 extern void dtrace_ret();

700 /* KPTI trampolines */
701 extern void tr_invaltrap();
702 extern void tr_div0trap(), tr_dbgtrap(), tr_nmiint(), tr_brktrap();
703 extern void tr_ovflotrap(), tr_boundstrap(), tr_invoptrap(), tr_ndptrap();
704 #if !defined(__xpv)
705 extern void tr_syserrtrap();
706 #endif

```

```

707 extern void tr_invaltrap(), tr_invtsstrap(), tr_segnptrap(), tr_stktrap();
708 extern void tr_gptrap(), tr_pfttrap(), tr_ndperr();
709 extern void tr_overrun(), tr_resvtrap();
710 extern void tr_achktrap(), tr_mcetrap();
711 extern void tr_xmtrap();
712 extern void tr_fasttrap();
713 extern void tr_dtrace_ret();

715 #if !defined(__amd64)
716 extern void pentium_pfttrap();
717 #endif

719 extern uint64_t kpti_enable;

721 #endif /* _ASM */

723 #ifdef __cplusplus
724 }

```

unchanged\_portion\_omitted

```

*****
32555 Thu Jun 14 17:16:00 2018
new/usr/src/uts/intel/sys/x86_archext.h
9600 LDT still not happy under KPTI
*****
_____unchanged_portion_omitted_____

731 extern int x86_use_pcid;
732 extern int x86_use_invpcid;

734 /*
735  * Utility functions to get/set extended control registers (XCR)
736  * Initial use is to get/set the contents of the XFEATURE_ENABLED_MASK.
737  */
738 extern uint64_t get_xcr(uint_t);
739 extern void set_xcr(uint_t, uint64_t);

741 extern uint64_t rdmsr(uint_t);
742 extern void wrmsr(uint_t, const uint64_t);
743 extern uint64_t xrdmsr(uint_t);
744 extern void xwrmsr(uint_t, const uint64_t);
745 extern int checked_rdmsr(uint_t, uint64_t *);
746 extern int checked_wrmsr(uint_t, uint64_t);

748 extern void invalidate_cache(void);
749 extern ulong_t getcr4(void);
750 extern void setcr4(ulong_t);

752 extern void mtrr_sync(void);

754 extern void cpu_fast_syscall_enable(void);
755 extern void cpu_fast_syscall_disable(void);
754 extern void cpu_fast_syscall_enable(void *);
755 extern void cpu_fast_syscall_disable(void *);

757 struct cpu;

759 extern int cpuid_checkpass(struct cpu *, int);
760 extern uint32_t cpuid_insn(struct cpu *, struct cpuid_regs *);
761 extern uint32_t __cpuid_insn(struct cpuid_regs *);
762 extern int cpuid_getbrandstr(struct cpu *, char *, size_t);
763 extern int cpuid_getidstr(struct cpu *, char *, size_t);
764 extern const char *cpuid_getvendorstr(struct cpu *);
765 extern uint_t cpuid_getvendor(struct cpu *);
766 extern uint_t cpuid_getfamily(struct cpu *);
767 extern uint_t cpuid_getmodel(struct cpu *);
768 extern uint_t cpuid_getstep(struct cpu *);
769 extern uint_t cpuid_getsig(struct cpu *);
770 extern uint_t cpuid_get_ncpu_per_chip(struct cpu *);
771 extern uint_t cpuid_get_ncore_per_chip(struct cpu *);
772 extern uint_t cpuid_get_ncpu_sharing_last_cache(struct cpu *);
773 extern id_t cpuid_get_last_lvl_cacheid(struct cpu *);
774 extern int cpuid_get_chipid(struct cpu *);
775 extern id_t cpuid_get_coreid(struct cpu *);
776 extern int cpuid_get_pkgcoreid(struct cpu *);
777 extern int cpuid_get_clogid(struct cpu *);
778 extern int cpuid_get_cacheid(struct cpu *);
779 extern uint32_t cpuid_get_apicid(struct cpu *);
780 extern uint_t cpuid_get_procnodid(struct cpu *cpu);
781 extern uint_t cpuid_get_procnodes_per_pkg(struct cpu *cpu);
782 extern uint_t cpuid_get_compunitid(struct cpu *cpu);
783 extern uint_t cpuid_get_cores_per_compunit(struct cpu *cpu);
784 extern size_t cpuid_get_xsave_size();
785 extern boolean_t cpuid_need_fp_excp_handling();
786 extern int cpuid_is_cmt(struct cpu *);
787 extern int cpuid_syscall132_insn(struct cpu *);

```

```

788 extern int getl2cacheinfo(struct cpu *, int *, int *, int *);

790 extern uint32_t cpuid_getchiprev(struct cpu *);
791 extern const char *cpuid_getchiprevstr(struct cpu *);
792 extern uint32_t cpuid_getsockettype(struct cpu *);
793 extern const char *cpuid_getsocketstr(struct cpu *);

795 extern int cpuid_have_cr8access(struct cpu *);

797 extern int cpuid_opteron_erratum(struct cpu *, uint_t);

799 struct cpuid_info;

801 extern void setx86isalist(void);
802 extern void cpuid_alloc_space(struct cpu *);
803 extern void cpuid_free_space(struct cpu *);
804 extern void cpuid_pass1(struct cpu *, uchar_t *);
805 extern void cpuid_pass2(struct cpu *);
806 extern void cpuid_pass3(struct cpu *);
807 extern void cpuid_pass4(struct cpu *, uint_t *);
808 extern void cpuid_set_cpu_properties(void *, processorid_t,
809     struct cpuid_info *);

811 extern void cpuid_get_addrsize(struct cpu *, uint_t *, uint_t *);
812 extern uint_t cpuid_get_dtlb_nent(struct cpu *, size_t);

814 #if !defined(__xpv)
815 extern uint32_t *cpuid_mwait_alloc(struct cpu *);
816 extern void cpuid_mwait_free(struct cpu *);
817 extern int cpuid_deep_cstates_supported(void);
818 extern int cpuid_arat_supported(void);
819 extern int cpuid_iepb_supported(struct cpu *);
820 extern int cpuid_deadline_tsc_supported(void);
821 extern void vmware_port(int, uint32_t *);
822 #endif

824 struct cpu_ucose_info;

826 extern void ucode_alloc_space(struct cpu *);
827 extern void ucode_free_space(struct cpu *);
828 extern void ucode_check(struct cpu *);
829 extern void ucode_cleanup();

831 #if !defined(__xpv)
832 extern char _tsc_mfence_start;
833 extern char _tsc_mfence_end;
834 extern char _tscp_start;
835 extern char _tscp_end;
836 extern char _no_rdtsc_start;
837 extern char _no_rdtsc_end;
838 extern char _tsc_lfence_start;
839 extern char _tsc_lfence_end;
840 #endif

842 #if !defined(__xpv)
843 extern char bcopy_patch_start;
844 extern char bcopy_patch_end;
845 extern char bcopy_ck_size;
846 #endif

848 extern void post_startup_cpu_fixups(void);

850 extern uint_t workaround_errata(struct cpu *);

852 #if defined(OPTERON_ERRATUM_93)
853 extern int opteron_erratum_93;

```

```
854 #endif

856 #if defined(OPTERON_ERRATUM_91)
857 extern int opteron_erratum_91;
858 #endif

860 #if defined(OPTERON_ERRATUM_100)
861 extern int opteron_erratum_100;
862 #endif

864 #if defined(OPTERON_ERRATUM_121)
865 extern int opteron_erratum_121;
866 #endif

868 #if defined(OPTERON_WORKAROUND_6323525)
869 extern int opteron_workaround_6323525;
870 extern void patch_workaround_6323525(void);
871 #endif

873 #if !defined(__xpv)
874 extern void determine_platform(void);
875 #endif
876 extern int get_hwenv(void);
877 extern int is_controldom(void);

879 extern void enable_pcid(void);

881 extern void xsave_setup_msr(struct cpu *);

883 /*
884  * Hypervisor signatures
885  */
886 #define HVSIG_XEN_HVM      "XenVMMXenVMM"
887 #define HVSIG_VMWARE      "VMwareVMware"
888 #define HVSIG_KVM         "KVMKVMKVM"
889 #define HVSIG_MICROSOFT   "Microsoft Hv"

891 /*
892  * Defined hardware environments
893  */
894 #define HW_NATIVE          (1 << 0)      /* Running on bare metal */
895 #define HW_XEN_PV         (1 << 1)      /* Running on Xen PVM */

897 #define HW_XEN_HVM        (1 << 2)      /* Running on Xen HVM */
898 #define HW_VMWARE        (1 << 3)      /* Running on VMware hypervisor */
899 #define HW_KVM           (1 << 4)      /* Running on KVM hypervisor */
900 #define HW_MICROSOFT     (1 << 5)      /* Running on Microsoft hypervisor */

902 #define HW_VIRTUAL        (HW_XEN_HVM | HW_VMWARE | HW_KVM | HW_MICROSOFT)

904 #endif /* _KERNEL */

906 #endif /* !_ASM */

908 /*
909  * VMware hypervisor related defines
910  */
911 #define VMWARE_HVMAGIC    0x564d5868
912 #define VMWARE_HVPORT    0x5658
913 #define VMWARE_HVCMD_GETVERSION 0x0a
914 #define VMWARE_HVCMD_GETTSCFREQ 0x2d

916 #ifdef __cplusplus
917 }

```

unchanged portion omitted