

new/usr/src/uts/common/cpr/cpr_dump.c

1

29374 Tue Nov 24 09:34:44 2015

new/usr/src/uts/common/cpr/cpr_dump.c

patch lower-case-segops

unchanged portion omitted

```
661 /*
662  * Count pages within each kernel segment; call cpr_sparse_seg_check()
663  * to find out whether a sparsely filled segment needs special
664  * treatment (e.g. kvseg).
665  * Todo: A "segop_cpr" like segop_dump should be introduced, the cpr
665  * Todo: A "SEGOP_CPR" like SEGOP_DUMP should be introduced, the cpr
666  * module shouldn't need to know segment details like if it is
667  * sparsely filled or not (makes kseg_table obsolete).
668  */
669 pgcnt_t
670 cpr_count_seg_pages(int mapflag, bitfunc_t bitfunc)
671 {
672     struct seg *segp;
673     pgcnt_t pages;
674     ksegtbl_entry_t *ste;
675
676     pages = 0;
677     for (segp = AS_SEGFIRST(&kas); segp; segp = AS_SEGNEXT(&kas, segp)) {
678         if (ste = cpr_sparse_seg_check(segp)) {
679             pages += (ste->st_fcn)(mapflag, bitfunc, segp);
680         } else {
681             pages += cpr_count_pages(segp->s_base,
682                                     segp->s_size, mapflag, bitfunc, DBG_SHOWRANGE);
683         }
684     }
685
686     return (pages);
687 }
```

unchanged portion omitted


```

1548     v[i].sh_offset = *doffsetp;
1549     v[i].sh_size = shdr->sh_size;
1550     if (symtab == NULL) {
1551         v[i].sh_link = 0;
1552     } else if (symtab->sh_type ==
1553               SHT_SYMTAB &&
1554               symtab_ndx != 0) {
1555         v[i].sh_link =
1556             symtab_ndx;
1557     } else {
1558         v[i].sh_link = i + 1;
1559     }
1561     copy_scn(shdr,.mvp, &v[i], vp,
1562             doffsetp, data, datasz, credp,
1563             rlimit);
1564 }
1566     ctf_ndx = i++;
1568     /*
1569     * We've already dumped the symtab.
1570     */
1571     if (symtab != NULL &&
1572         symtab->sh_type == SHT_SYMTAB &&
1573         symtab_ndx != 0)
1574         continue;
1576     } else if (strcmp(name,
1577                   shstrtab_data[STR_SYMTAB]) == 0) {
1578         if ((content & CC_CONTENT_SYMTAB) == 0 ||
1579             symtab != 0)
1580             continue;
1582         symtab = shdr;
1583     }
1585     if (symtab != NULL) {
1586         if ((symtab->sh_type != SHT_DYNSYM &&
1587             symtab->sh_type != SHT_SYMTAB) ||
1588             symtab->sh_link == 0 ||
1589             symtab->sh_link >= nshdrs)
1590             continue;
1592         strtab = (Shdr *) (shbase +
1593                          symtab->sh_link * ehdr.e_shentsize);
1595         if (strtab->sh_type != SHT_STRTAB)
1596             continue;
1598         if (v != NULL && i < nv - 2) {
1599             sz = MAX(symtab->sh_size,
1600                   strtab->sh_size);
1601             if (sz > datasz &&
1602                 sz <= elf_datasz_max) {
1603                 if (data != NULL)
1604                     kmem_free(data, datasz);
1606                 datasz = sz;
1607                 data = kmem_alloc(datasz,
1608                                 KM_SLEEP);
1609             }
1611             if (symtab->sh_type == SHT_DYNSYM) {
1612                 v[i].sh_name = shstrtab_ndx(
1613                     &shstrtab, STR_DYNSYM);

```

```

1614         v[i + 1].sh_name = shstrtab_ndx(
1615             &shstrtab, STR_DYNSTR);
1616     } else {
1617         v[i].sh_name = shstrtab_ndx(
1618             &shstrtab, STR_SYMTAB);
1619         v[i + 1].sh_name = shstrtab_ndx(
1620             &shstrtab, STR_STRTAB);
1621     }
1623     v[i].sh_type = symtab->sh_type;
1624     v[i].sh_addr = symtab->sh_addr;
1625     if (ehdr.e_type == ET_DYN ||
1626         v[i].sh_addr == 0)
1627         v[i].sh_addr +=
1628             (Addr)(uintptr_t)saddr;
1629     v[i].sh_addralign =
1630         symtab->sh_addralign;
1631     *doffsetp = roundup(*doffsetp,
1632                       v[i].sh_addralign);
1633     v[i].sh_offset = *doffsetp;
1634     v[i].sh_size = symtab->sh_size;
1635     v[i].sh_link = i + 1;
1636     v[i].sh_entsize = symtab->sh_entsize;
1637     v[i].sh_info = symtab->sh_info;
1639     copy_scn(symtab,.mvp, &v[i], vp,
1640             doffsetp, data, datasz, credp,
1641             rlimit);
1643     v[i + 1].sh_type = SHT_STRTAB;
1644     v[i + 1].sh_flags = SHF_STRINGS;
1645     v[i + 1].sh_addr = symtab->sh_addr;
1646     if (ehdr.e_type == ET_DYN ||
1647         v[i + 1].sh_addr == 0)
1648         v[i + 1].sh_addr +=
1649             (Addr)(uintptr_t)saddr;
1650     v[i + 1].sh_addralign =
1651         strtab->sh_addralign;
1652     *doffsetp = roundup(*doffsetp,
1653                       v[i + 1].sh_addralign);
1654     v[i + 1].sh_offset = *doffsetp;
1655     v[i + 1].sh_size = strtab->sh_size;
1657     copy_scn(strtab,.mvp, &v[i + 1], vp,
1658             doffsetp, data, datasz, credp,
1659             rlimit);
1660 }
1662     if (symtab->sh_type == SHT_SYMTAB)
1663         symtab_ndx = i;
1664         i += 2;
1665     }
1666 }
1668     kmem_free(shstrbase, shstrsize);
1669     kmem_free(shbase, shsize);
1671     lastvp =.mvp;
1672 }
1674     if (v == NULL) {
1675         if (i == 1)
1676             *nshdrsp = 0;
1677         else
1678             *nshdrsp = i + 1;
1679         goto done;

```

```

1680     }
1682     if (i != nv - 1) {
1683         cmn_err(CE_WARN, "elfcore: core dump failed for "
1684             "process %d; address space is changing", p->p_pid);
1685         error = EIO;
1686         goto done;
1687     }
1689     v[i].sh_name = shstrtab_ndx(&shstrtab, STR_SHSTRTAB);
1690     v[i].sh_size = shstrtab_size(&shstrtab);
1691     v[i].sh_addralign = 1;
1692     *doffsetp = roundup(*doffsetp, v[i].sh_addralign);
1693     v[i].sh_offset = *doffsetp;
1694     v[i].sh_flags = SHF_STRINGS;
1695     v[i].sh_type = SHT_STRTAB;
1697     if (v[i].sh_size > datasz) {
1698         if (data != NULL)
1699             kmem_free(data, datasz);
1701         datasz = v[i].sh_size;
1702         data = kmem_alloc(datasz,
1703             KM_SLEEP);
1704     }
1706     shstrtab_dump(&shstrtab, data);
1708     if ((error = core_write(vp, UIO_SYSSPACE, *doffsetp,
1709         data, v[i].sh_size, rlimit, credp)) != 0)
1710         goto done;
1712     *doffsetp += v[i].sh_size;
1714 done:
1715     if (data != NULL)
1716         kmem_free(data, datasz);
1718     return (error);
1719 }
1721 int
1722 elfcore(vnode_t *vp, proc_t *p, cred_t *credp, rlim64_t rlimit, int sig,
1723     core_content_t content)
1724 {
1725     offset_t poffset, soffset;
1726     Off doffset;
1727     int error, i, nphdrs, nshdrs;
1728     int overflow = 0;
1729     struct seg *seg;
1730     struct as *as = p->p_as;
1731     union {
1732         Ehdr ehdr;
1733         Phdr phdr[1];
1734         Shdr shdr[1];
1735     } *bigwad;
1736     size_t bigsize;
1737     size_t phdrsz, shdrsz;
1738     Ehdr *ehdr;
1739     Phdr *v;
1740     caddr_t brkbase;
1741     size_t brksize;
1742     caddr_t stkbase;
1743     size_t stksize;
1744     int ntries = 0;
1745     klwp_t *lwp = ttolwp(curthread);

```

```

1747 top:
1748     /*
1749     * Make sure we have everything we need (registers, etc.).
1750     * All other lwps have already stopped and are in an orderly state.
1751     */
1752     ASSERT(p == ttoproc(curthread));
1753     prstop(0, 0);
1755     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
1756     nphdrs = prnsegs(as, 0) + 2; /* two CORE note sections */
1758     /*
1759     * Count the number of section headers we're going to need.
1760     */
1761     nshdrs = 0;
1762     if (content & (CC_CONTENT_CTF | CC_CONTENT_SYMTAB)) {
1763         (void) process_scns(content, p, credp, NULL, NULL, NULL, 0,
1764             NULL, &nshdrs);
1765     }
1766     AS_LOCK_EXIT(as, &as->a_lock);
1768     ASSERT(nshdrs == 0 || nshdrs > 1);
1770     /*
1771     * The core file contents may required zero section headers, but if
1772     * we overflow the 16 bits allotted to the program header count in
1773     * the ELF header, we'll need that program header at index zero.
1774     */
1775     if (nshdrs == 0 && nphdrs >= PN_XNUM)
1776         nshdrs = 1;
1778     phdrsz = nphdrs * sizeof (Phdr);
1779     shdrsz = nshdrs * sizeof (Shdr);
1781     bigsize = MAX(sizeof (*bigwad), MAX(phdrsz, shdrsz));
1782     bigwad = kmem_alloc(bigsize, KM_SLEEP);
1784     ehdr = &bigwad->ehdr;
1785     bzero(ehdr, sizeof (*ehdr));
1787     ehdr->e_ident[EI_MAG0] = ELFMAG0;
1788     ehdr->e_ident[EI_MAG1] = ELFMAG1;
1789     ehdr->e_ident[EI_MAG2] = ELFMAG2;
1790     ehdr->e_ident[EI_MAG3] = ELFMAG3;
1791     ehdr->e_ident[EI_CLASS] = ELFCLASS;
1792     ehdr->e_type = ET_CORE;
1794 #if !defined(_LP64) || defined(_ELF32_COMPAT)
1796 #if defined(__sparc)
1797     ehdr->e_ident[EI_DATA] = ELFDATA2MSB;
1798     ehdr->e_machine = EM_SPARC;
1799 #elif defined(__i386) || defined(__i386_COMPAT)
1800     ehdr->e_ident[EI_DATA] = ELFDATA2LSB;
1801     ehdr->e_machine = EM_386;
1802 #else
1803 #error "no recognized machine type is defined"
1804 #endif
1806 #else /* !defined(_LP64) || defined(_ELF32_COMPAT) */
1808 #if defined(__sparc)
1809     ehdr->e_ident[EI_DATA] = ELFDATA2MSB;
1810     ehdr->e_machine = EM_SPARCV9;
1811 #elif defined(__amd64)

```

```

1812     ehdr->e_ident[EI_DATA] = ELFDATA2LSB;
1813     ehdr->e_machine = EM_AMD64;
1814 #else
1815 #error "no recognized 64-bit machine type is defined"
1816 #endif

1818 #endif /* !defined(_LP64) || defined(_ELF32_COMPAT) */

1820 /*
1821  * If the count of program headers or section headers or the index
1822  * of the section string table can't fit in the mere 16 bits
1823  * shortsightedly allotted to them in the ELF header, we use the
1824  * extended formats and put the real values in the section header
1825  * as index 0.
1826  */
1827 ehdr->e_version = EV_CURRENT;
1828 ehdr->e_ehsize = sizeof (Ehdr);

1830 if (nphdrs >= PN_XNUM)
1831     ehdr->e_phnum = PN_XNUM;
1832 else
1833     ehdr->e_phnum = (unsigned short)nphdrs;

1835 ehdr->e_phoff = sizeof (Ehdr);
1836 ehdr->e_phentsize = sizeof (Phdr);

1838 if (nshdrs > 0) {
1839     if (nshdrs >= SHN_LORESERVE)
1840         ehdr->e_shnum = 0;
1841     else
1842         ehdr->e_shnum = (unsigned short)nshdrs;

1844     if (nshdrs - 1 >= SHN_LORESERVE)
1845         ehdr->e_shstrndx = SHN_XINDEX;
1846     else
1847         ehdr->e_shstrndx = (unsigned short)(nshdrs - 1);

1849     ehdr->e_shoff = ehdr->e_phoff + ehdr->e_phentsize * nphdrs;
1850     ehdr->e_shentsize = sizeof (Shdr);
1851 }

1853 if (error = core_write(vp, UIO_SYSSPACE, (offset_t)0, ehdr,
1854     sizeof (Ehdr), rlimit, credp))
1855     goto done;

1857 poffset = sizeof (Ehdr);
1858 soffset = sizeof (Ehdr) + phdrsz;
1859 doffset = sizeof (Ehdr) + phdrsz + shdrsz;

1861 v = &bigwad->phdr[0];
1862 bzero(v, phdrsz);

1864 setup_old_note_header(&v[0], p);
1865 v[0].p_offset = doffset = roundup(doffset, sizeof (Word));
1866 doffset += v[0].p_filesz;

1868 setup_note_header(&v[1], p);
1869 v[1].p_offset = doffset = roundup(doffset, sizeof (Word));
1870 doffset += v[1].p_filesz;

1872 mutex_enter(&p->p_lock);

1874 brkbase = p->p_brkbase;
1875 brksize = p->p_brksize;

1877 stkbase = p->p_usrstack - p->p_stksize;

```

```

1878     stksize = p->p_stksize;

1880     mutex_exit(&p->p_lock);

1882     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
1883     i = 2;
1884     for (seg = AS_SEGFIRST(as); seg != NULL; seg = AS_SEGNEXT(as, seg)) {
1885         caddr_t eaddr = seg->s_base + pr_getsegsz(seg, 0);
1886         caddr_t saddr, naddr;
1887         void *tmp = NULL;
1888         extern struct seg_ops segspt_shmops;

1890         for (saddr = seg->s_base; saddr < eaddr; saddr = naddr) {
1891             uint_t prot;
1892             size_t size;
1893             int type;
1894             vnode_t *mvp;

1896             prot = pr_getprot(seg, 0, &tmp, &saddr, &naddr, eaddr);
1897             prot &= PROT_READ | PROT_WRITE | PROT_EXEC;
1898             if ((size = (size_t)(naddr - saddr)) == 0)
1899                 continue;
1900             if (i == nphdrs) {
1901                 overflow++;
1902                 continue;
1903             }
1904             v[i].p_type = PT_LOAD;
1905             v[i].p_vaddr = (Addr)(uintptr_t)saddr;
1906             v[i].p_memsz = size;
1907             if (prot & PROT_READ)
1908                 v[i].p_flags |= PF_R;
1909             if (prot & PROT_WRITE)
1910                 v[i].p_flags |= PF_W;
1911             if (prot & PROT_EXEC)
1912                 v[i].p_flags |= PF_X;

1914             /*
1915              * Figure out which mappings to include in the core.
1916              */
1917             type = segop_gettype(seg, saddr);
1917             type = SEGOP_GETTYPE(seg, saddr);

1919             if (saddr == stkbase && size == stksize) {
1920                 if (!(content & CC_CONTENT_STACK))
1921                     goto exclude;

1923             } else if (saddr == brkbase && size == brksize) {
1924                 if (!(content & CC_CONTENT_HEAP))
1925                     goto exclude;

1927             } else if (seg->s_ops == &segspt_shmops) {
1928                 if (type & MAP_NORESERVE) {
1929                     if (!(content & CC_CONTENT_DISM))
1930                         goto exclude;
1931                 } else {
1932                     if (!(content & CC_CONTENT_ISM))
1933                         goto exclude;
1934                 }

1936             } else if (seg->s_ops != &segvn_ops) {
1937                 goto exclude;

1939             } else if (type & MAP_SHARED) {
1940                 if (shmgetid(p, saddr) != SHMID_NONE) {
1941                     if (!(content & CC_CONTENT_SHM))
1942                         goto exclude;

```

```

1944         } else if (segop_getvp(seg, seg->s_base,
1944         ) else if (SEGOP_GETVP(seg, seg->s_base,
1945         &mvp) != 0 || mvp == NULL ||
1946         mvp->v_type != VREG) {
1947             if (!(content & CC_CONTENT_SHANON))
1948                 goto exclude;
1949
1950         } else {
1951             if (!(content & CC_CONTENT_SHFILE))
1952                 goto exclude;
1953         }
1954
1955         } else if (segop_getvp(seg, seg->s_base, &mvp) != 0 ||
1955         ) else if (SEGOP_GETVP(seg, seg->s_base, &mvp) != 0 ||
1956         mvp == NULL || mvp->v_type != VREG) {
1957             if (!(content & CC_CONTENT_ANON))
1958                 goto exclude;
1959
1960         } else if (prot == (PROT_READ | PROT_EXEC)) {
1961             if (!(content & CC_CONTENT_TEXT))
1962                 goto exclude;
1963
1964         } else if (prot == PROT_READ) {
1965             if (!(content & CC_CONTENT_RODATA))
1966                 goto exclude;
1967
1968         } else {
1969             if (!(content & CC_CONTENT_DATA))
1970                 goto exclude;
1971         }
1972
1973         doffset = roundup(doffset, sizeof (Word));
1974         v[i].p_offset = doffset;
1975         v[i].p_filesz = size;
1976         doffset += size;
1977     exclude:
1978         i++;
1979     }
1980     ASSERT(tmp == NULL);
1981 }
1982 AS_LOCK_EXIT(as, &as->a_lock);
1983
1984 if (overflow || i != nphdrs) {
1985     if (ntries++ == 0) {
1986         kmem_free(bigwad, bigsize);
1987         overflow = 0;
1988         goto top;
1989     }
1990     cmn_err(CE_WARN, "elfcore: core dump failed for "
1991            "process %d; address space is changing", p->p_pid);
1992     error = EIO;
1993     goto done;
1994 }
1995
1996 if ((error = core_write(vp, UIO_SYSSPACE, poffset,
1997     v, phdrsz, rlimit, credp)) != 0)
1998     goto done;
1999
2000 if ((error = write_old_elfnotes(p, sig, vp, v[0].p_offset, rlimit,
2001     credp)) != 0)
2002     goto done;
2003
2004 if ((error = write_elfnotes(p, sig, vp, v[1].p_offset, rlimit,
2005     credp, content)) != 0)
2006     goto done;

```

```

2008         for (i = 2; i < nphdrs; i++) {
2009             prkillinfo_t killinfo;
2010             sigqueue_t *sq;
2011             int sig, j;
2012
2013             if (v[i].p_filesz == 0)
2014                 continue;
2015
2016             /*
2017             * If dumping out this segment fails, rather than failing
2018             * the core dump entirely, we reset the size of the mapping
2019             * to zero to indicate that the data is absent from the core
2020             * file and or in the PF_SUNW_FAILURE flag to differentiate
2021             * this from mappings that were excluded due to the core file
2022             * content settings.
2023             */
2024             if ((error = core_seg(p, vp, v[i].p_offset,
2025                 (caddr_t)(uintptr_t)v[i].p_vaddr, v[i].p_filesz,
2026                 rlimit, credp)) == 0) {
2027                 continue;
2028             }
2029
2030             if ((sig = lwp->lwp_cursig) == 0) {
2031                 /*
2032                 * We failed due to something other than a signal.
2033                 * Since the space reserved for the segment is now
2034                 * unused, we stash the errno in the first four
2035                 * bytes. This undocumented interface will let us
2036                 * understand the nature of the failure.
2037                 */
2038                 (void) core_write(vp, UIO_SYSSPACE, v[i].p_offset,
2039                     &error, sizeof (error), rlimit, credp);
2040
2041                 v[i].p_filesz = 0;
2042                 v[i].p_flags |= PF_SUNW_FAILURE;
2043                 if ((error = core_write(vp, UIO_SYSSPACE,
2044                     poffset + sizeof (v[i]) * i, &v[i], sizeof (v[i]),
2045                     rlimit, credp)) != 0)
2046                     goto done;
2047
2048                 continue;
2049             }
2050
2051             /*
2052             * We took a signal. We want to abort the dump entirely, but
2053             * we also want to indicate what failed and why. We therefore
2054             * use the space reserved for the first failing segment to
2055             * write our error (which, for purposes of compatability with
2056             * older core dump readers, we set to EINTR) followed by any
2057             * siginfo associated with the signal.
2058             */
2059             bzero(&killinfo, sizeof (killinfo));
2060             killinfo.prk_error = EINTR;
2061
2062             sq = sig == SIGKILL ? curproc->p_killsq : lwp->lwp_curinfo;
2063
2064             if (sq != NULL) {
2065                 bcopy(&sq->sq_info, &killinfo.prk_info,
2066                     sizeof (sq->sq_info));
2067             } else {
2068                 killinfo.prk_info.si_signo = lwp->lwp_cursig;
2069                 killinfo.prk_info.si_code = SI_NOINFO;
2070             }
2071
2072 #if (defined(_SYSCALL32_IMPL) || defined(_LP64))

```

```

2073      /*
2074      * If this is a 32-bit process, we need to translate from the
2075      * native siginfo to the 32-bit variant. (Core readers must
2076      * always have the same data model as their target or must
2077      * be aware of -- and compensate for -- data model differences.)
2078      */
2079      if (curproc->p_model == DATAMODEL_ILP32) {
2080          siginfo32_t si32;
2081
2082          siginfo_kto32((k_siginfo_t *)&killinfo.prk_info, &si32);
2083          bcopy(&si32, &killinfo.prk_info, sizeof (si32));
2084      }
2085 #endif
2086
2087      (void) core_write(vp, UIO_SYSSPACE, v[i].p_offset,
2088                      &killinfo, sizeof (killinfo), rlimit, credp);
2089
2090      /*
2091      * For the segment on which we took the signal, indicate that
2092      * its data now refers to a siginfo.
2093      */
2094      v[i].p_filesz = 0;
2095      v[i].p_flags |= PF_SUNW_FAILURE | PF_SUNW_KILLED |
2096                  PF_SUNW_SIGINFO;
2097
2098      /*
2099      * And for every other segment, indicate that its absence
2100      * is due to a signal.
2101      */
2102      for (j = i + 1; j < nphdrs; j++) {
2103          v[j].p_filesz = 0;
2104          v[j].p_flags |= PF_SUNW_FAILURE | PF_SUNW_KILLED;
2105      }
2106
2107      /*
2108      * Finally, write out our modified program headers.
2109      */
2110      if ((error = core_write(vp, UIO_SYSSPACE,
2111                             poffset + sizeof (v[i]) * i, &v[i],
2112                             sizeof (v[i]) * (nphdrs - i), rlimit, credp)) != 0)
2113          goto done;
2114
2115      break;
2116  }
2117
2118  if (nshdrs > 0) {
2119      bzero(&bigwad->shdr[0], shdrsz);
2120
2121      if (nshdrs >= SHN_LORESERVE)
2122          bigwad->shdr[0].sh_size = nshdrs;
2123
2124      if (nshdrs - 1 >= SHN_LORESERVE)
2125          bigwad->shdr[0].sh_link = nshdrs - 1;
2126
2127      if (nphdrs >= PN_XNUM)
2128          bigwad->shdr[0].sh_info = nphdrs;
2129
2130      if (nshdrs > 1) {
2131          AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
2132          if ((error = process_scns(content, p, credp, vp,
2133                                  &bigwad->shdr[0], nshdrs, rlimit, &doffset,
2134                                  NULL)) != 0) {
2135              AS_LOCK_EXIT(as, &as->a_lock);
2136              goto done;
2137          }
2138          AS_LOCK_EXIT(as, &as->a_lock);

```

```

2139      }
2140
2141      if ((error = core_write(vp, UIO_SYSSPACE, soffset,
2142                             &bigwad->shdr[0], shdrsz, rlimit, credp)) != 0)
2143          goto done;
2144      }
2145
2146  done:
2147      kmem_free(bigwad, bigsize);
2148      return (error);
2149  }

```

unchanged portion omitted

```

*****
171811 Tue Nov 24 09:34:44 2015
new/usr/src/uts/common/fs/nfs/nfs3_vnops.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

5543 /*
5544  * Setup and add an address space callback to do the work of the delmap call.
5545  * The callback will (and must be) deleted in the actual callback function.
5546  *
5547  * This is done in order to take care of the problem that we have with holding
5548  * the address space's a_lock for a long period of time (e.g. if the NFS server
5549  * is down). Callbacks will be executed in the address space code while the
5550  * a_lock is not held. Holding the address space's a_lock causes things such
5551  * as ps and fork to hang because they are trying to acquire this lock as well.
5552  */
5553 /* ARGSUSED */
5554 static int
5555 nfs3_delmap(vnode_t *vp, offset_t off, struct as *as, caddr_t addr,
5556             size_t len, uint_t prot, uint_t maxprot, uint_t flags,
5557             cred_t *cr, caller_context_t *ct)
5558 {
5559     int                caller_found;
5560     int                error;
5561     rnode_t           *rp;
5562     nfs_delmap_args_t *dmapp;
5563     nfs_delmapcall_t  *delmap_call;

5565     if (vp->v_flag & VNOMAP)
5566         return (ENOSYS);
5567     /*
5568      * A process may not change zones if it has NFS pages mmap'ed
5569      * in, so we can't legitimately get here from the wrong zone.
5570      */
5571     ASSERT(nfs_zone() == VTOMI(vp)->mi_zone);

5573     rp = VTOR(vp);

5575     /*
5576      * The way that the address space of this process deletes its mapping
5577      * of this file is via the following call chains:
5578      * - as_free()->segop_unmap()/segvn_unmap()->VOP_DELMAP()/nfs3_delmap()
5579      * - as_unmap()->segop_unmap()/segvn_unmap()->VOP_DELMAP()/nfs3_delmap()
5580      * - as_free()->SEGOP_UNMAP()/segvn_unmap()->VOP_DELMAP()/nfs3_delmap()
5581      * - as_unmap()->SEGOP_UNMAP()/segvn_unmap()->VOP_DELMAP()/nfs3_delmap()
5582      *
5583      * With the use of address space callbacks we are allowed to drop the
5584      * address space lock, a_lock, while executing the NFS operations that
5585      * need to go over the wire. Returning EAGAIN to the caller of this
5586      * function is what drives the execution of the callback that we add
5587      * below. The callback will be executed by the address space code
5588      * after dropping the a_lock. When the callback is finished, since
5589      * we dropped the a_lock, it must be re-acquired and segvn_unmap()
5590      * is called again on the same segment to finish the rest of the work
5591      * that needs to happen during unmapping.
5592      *
5593      * This action of calling back into the segment driver causes
5594      * nfs3_delmap() to get called again, but since the callback was
5595      * already executed at this point, it already did the work and there
5596      * is nothing left for us to do.
5597      *
5598      * To Summarize:
5599      * - The first time nfs3_delmap is called by the current thread is when
5600      * we add the caller associated with this delmap to the delmap caller
5601      * list, add the callback, and return EAGAIN.

```

```

5600     * - The second time in this call chain when nfs3_delmap is called we
5601     * will find this caller in the delmap caller list and realize there
5602     * is no more work to do thus removing this caller from the list and
5603     * returning the error that was set in the callback execution.
5604     */
5605     caller_found = nfs_find_and_delete_delmapcall(rp, &error);
5606     if (caller_found) {
5607         /*
5608          * 'error' is from the actual delmap operations. To avoid
5609          * hangs, we need to handle the return of EAGAIN differently
5610          * since this is what drives the callback execution.
5611          * In this case, we don't want to return EAGAIN and do the
5612          * callback execution because there are none to execute.
5613          */
5614         if (error == EAGAIN)
5615             return (0);
5616         else
5617             return (error);
5618     }

5620     /* current caller was not in the list */
5621     delmap_call = nfs_init_delmapcall();

5623     mutex_enter(&rp->r_statelock);
5624     list_insert_tail(&rp->r_indelmap, delmap_call);
5625     mutex_exit(&rp->r_statelock);

5627     dmapp = kmem_alloc(sizeof (nfs_delmap_args_t), KM_SLEEP);

5629     dmapp->vp = vp;
5630     dmapp->off = off;
5631     dmapp->addr = addr;
5632     dmapp->len = len;
5633     dmapp->prot = prot;
5634     dmapp->maxprot = maxprot;
5635     dmapp->flags = flags;
5636     dmapp->cr = cr;
5637     dmapp->caller = delmap_call;

5639     error = as_add_callback(as, nfs3_delmap_callback, dmapp,
5640                           AS_UNMAP_EVENT, addr, len, KM_SLEEP);

5642     return (error ? error : EAGAIN);
5643 }
_____unchanged_portion_omitted_____

```



```

*****
429800 Tue Nov 24 09:34:44 2015
new/usr/src/uts/common/fs/nfs/nfs4_vnops.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

11026 /*
11027  * Setup and add an address space callback to do the work of the delmap call.
11028  * The callback will (and must be) deleted in the actual callback function.
11029  *
11030  * This is done in order to take care of the problem that we have with holding
11031  * the address space's a_lock for a long period of time (e.g. if the NFS server
11032  * is down). Callbacks will be executed in the address space code while the
11033  * a_lock is not held. Holding the address space's a_lock causes things such
11034  * as ps and fork to hang because they are trying to acquire this lock as well.
11035  */
11036 /* ARGSUSED */
11037 static int
11038 nfs4_delmap(vnode_t *vp, offset_t off, struct as *as, caddr_t addr,
11039             size_t len, uint_t prot, uint_t maxprot, uint_t flags, cred_t *cr,
11040             caller_context_t *ct)
11041 {
11042     int                caller_found;
11043     int                error;
11044     rnode4_t          *rp;
11045     nfs4_delmap_args_t *dmapp;
11046     nfs4_delmapcall_t *delmap_call;

11048     if (vp->v_flag & VNOMAP)
11049         return (ENOSYS);

11051     /*
11052     * A process may not change zones if it has NFS pages mmap'ed
11053     * in, so we can't legitimately get here from the wrong zone.
11054     */
11055     ASSERT(nfs_zone() == VTOMI4(vp)->mi_zone);

11057     rp = VTOR4(vp);

11059     /*
11060     * The way that the address space of this process deletes its mapping
11061     * of this file is via the following call chains:
11062     * - as_free()->segop_unmap()/segvn_unmap()->VOP_DELMAP()/nfs4_delmap()
11063     * - as_unmap()->segop_unmap()/segvn_unmap()->VOP_DELMAP()/nfs4_delmap()
11064     * - as_free()->SEGOP_UNMAP()/segvn_unmap()->VOP_DELMAP()/nfs4_delmap()
11065     * - as_unmap()->SEGOP_UNMAP()/segvn_unmap()->VOP_DELMAP()/nfs4_delmap()
11066     *
11067     * With the use of address space callbacks we are allowed to drop the
11068     * address space lock, a_lock, while executing the NFS operations that
11069     * need to go over the wire. Returning EAGAIN to the caller of this
11070     * function is what drives the execution of the callback that we add
11071     * below. The callback will be executed by the address space code
11072     * after dropping the a_lock. When the callback is finished, since
11073     * we dropped the a_lock, it must be re-acquired and segvn_unmap()
11074     * is called again on the same segment to finish the rest of the work
11075     * that needs to happen during unmapping.
11076     *
11077     * This action of calling back into the segment driver causes
11078     * nfs4_delmap() to get called again, but since the callback was
11079     * already executed at this point, it already did the work and there
11080     * is nothing left for us to do.
11081     *
11082     * To Summarize:
11083     * - The first time nfs4_delmap is called by the current thread is when
11084     * we add the caller associated with this delmap to the delmap caller

```

```

11083     * list, add the callback, and return EAGAIN.
11084     * - The second time in this call chain when nfs4_delmap is called we
11085     * will find this caller in the delmap caller list and realize there
11086     * is no more work to do thus removing this caller from the list and
11087     * returning the error that was set in the callback execution.
11088     */
11089     caller_found = nfs4_find_and_delete_delmapcall(rp, &error);
11090     if (caller_found) {
11091         /*
11092         * 'error' is from the actual delmap operations. To avoid
11093         * hangs, we need to handle the return of EAGAIN differently
11094         * since this is what drives the callback execution.
11095         * In this case, we don't want to return EAGAIN and do the
11096         * callback execution because there are none to execute.
11097         */
11098         if (error == EAGAIN)
11099             return (0);
11100         else
11101             return (error);
11102     }

11104     /* current caller was not in the list */
11105     delmap_call = nfs4_init_delmapcall();

11107     mutex_enter(&rp->r_statelock);
11108     list_insert_tail(&rp->r_indelemap, delmap_call);
11109     mutex_exit(&rp->r_statelock);

11111     dmapp = kmem_alloc(sizeof (nfs4_delmap_args_t), KM_SLEEP);

11113     dmapp->vp = vp;
11114     dmapp->off = off;
11115     dmapp->addr = addr;
11116     dmapp->len = len;
11117     dmapp->prot = prot;
11118     dmapp->maxprot = maxprot;
11119     dmapp->flags = flags;
11120     dmapp->cr = cr;
11121     dmapp->caller = delmap_call;

11123     error = as_add_callback(as, nfs4_delmap_callback, dmapp,
11124                          AS_UNMAP_EVENT, addr, len, KM_SLEEP);

11126     return (error ? error : EAGAIN);
11127 }
_____unchanged_portion_omitted_____

```

```

*****
130899 Tue Nov 24 09:34:45 2015
new/usr/src/uts/common/fs/nfs/nfs_vnops.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

4640 /*
4641  * Setup and add an address space callback to do the work of the delmap call.
4642  * The callback will (and must be) deleted in the actual callback function.
4643  *
4644  * This is done in order to take care of the problem that we have with holding
4645  * the address space's a_lock for a long period of time (e.g. if the NFS server
4646  * is down). Callbacks will be executed in the address space code while the
4647  * a_lock is not held. Holding the address space's a_lock causes things such
4648  * as ps and fork to hang because they are trying to acquire this lock as well.
4649  */
4650 /* ARGSUSED */
4651 static int
4652 nfs_delmap(vnode_t *vp, offset_t off, struct as *as, caddr_t addr,
4653            size_t len, uint_t prot, uint_t maxprot, uint_t flags, cred_t *cr,
4654            caller_context_t *ct)
4655 {
4656     int                caller_found;
4657     int                error;
4658     rnode_t            *rp;
4659     nfs_delmap_args_t  *dmapp;
4660     nfs_delmapcall_t   *delmap_call;

4662     if (vp->v_flag & VNOMAP)
4663         return (ENOSYS);
4664     /*
4665      * A process may not change zones if it has NFS pages mmap'ed
4666      * in, so we can't legitimately get here from the wrong zone.
4667      */
4668     ASSERT(nfs_zone() == VTOMI(vp)->mi_zone);

4670     rp = VTOR(vp);

4672     /*
4673      * The way that the address space of this process deletes its mapping
4674      * of this file is via the following call chains:
4675      * - as_free()->segop_unmap()/segvn_unmap()->VOP_DELMAP()/nfs_delmap()
4676      * - as_unmap()->segop_unmap()/segvn_unmap()->VOP_DELMAP()/nfs_delmap()
4677      * - as_free()->SEGOP_UNMAP()/segvn_unmap()->VOP_DELMAP()/nfs_delmap()
4678      * - as_unmap()->SEGOP_UNMAP()/segvn_unmap()->VOP_DELMAP()/nfs_delmap()
4679      *
4680      * With the use of address space callbacks we are allowed to drop the
4681      * address space lock, a_lock, while executing the NFS operations that
4682      * need to go over the wire. Returning EAGAIN to the caller of this
4683      * function is what drives the execution of the callback that we add
4684      * below. The callback will be executed by the address space code
4685      * after dropping the a_lock. When the callback is finished, since
4686      * we dropped the a_lock, it must be re-acquired and segvn_unmap()
4687      * is called again on the same segment to finish the rest of the work
4688      * that needs to happen during unmapping.
4689      *
4690      * This action of calling back into the segment driver causes
4691      * nfs_delmap() to get called again, but since the callback was
4692      * already executed at this point, it already did the work and there
4693      * is nothing left for us to do.
4694      *
4695      * To Summarize:
4696      * - The first time nfs_delmap is called by the current thread is when
4697      * we add the caller associated with this delmap to the delmap caller
4698      * list, add the callback, and return EAGAIN.

```

```

4697     * - The second time in this call chain when nfs_delmap is called we
4698     * will find this caller in the delmap caller list and realize there
4699     * is no more work to do thus removing this caller from the list and
4700     * returning the error that was set in the callback execution.
4701     */
4702     caller_found = nfs_find_and_delete_delmapcall(rp, &error);
4703     if (caller_found) {
4704         /*
4705          * 'error' is from the actual delmap operations. To avoid
4706          * hangs, we need to handle the return of EAGAIN differently
4707          * since this is what drives the callback execution.
4708          * In this case, we don't want to return EAGAIN and do the
4709          * callback execution because there are none to execute.
4710          */
4711         if (error == EAGAIN)
4712             return (0);
4713         else
4714             return (error);
4715     }

4717     /* current caller was not in the list */
4718     delmap_call = nfs_init_delmapcall();

4720     mutex_enter(&rp->r_statelock);
4721     list_insert_tail(&rp->r_indeimap, delmap_call);
4722     mutex_exit(&rp->r_statelock);

4724     dmapp = kmem_alloc(sizeof (nfs_delmap_args_t), KM_SLEEP);

4726     dmapp->vp = vp;
4727     dmapp->off = off;
4728     dmapp->addr = addr;
4729     dmapp->len = len;
4730     dmapp->prot = prot;
4731     dmapp->maxprot = maxprot;
4732     dmapp->flags = flags;
4733     dmapp->cr = cr;
4734     dmapp->caller = delmap_call;

4736     error = as_add_callback(as, nfs_delmap_callback, dmapp,
4737                           AS_UNMAP_EVENT, addr, len, KM_SLEEP);

4739     return (error ? error : EAGAIN);
4740 }
_____unchanged_portion_omitted_____

```

```

*****
93906 Tue Nov 24 09:34:45 2015
new/usr/src/uts/common/fs/proc/priocntl.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

3117 /*
3118  * Common code for PIOCOPENM
3119  * Returns with the process unlocked.
3120  */
3121 static int
3122 propenm(prnode_t *pnp, caddr_t cmaddr, caddr_t va, int *rvalp, cred_t *cr)
3123 {
3124     proc_t *p = pnp->pr_common->prc_proc;
3125     struct as *as = p->p_as;
3126     int error = 0;
3127     struct seg *seg;
3128     struct vnode *xvp;
3129     int n;

3131     /*
3132     * By fiat, a system process has no address space.
3133     */
3134     if ((p->p_flag & SSYS) || as == &kas) {
3135         error = EINVAL;
3136     } else if (cmaddr) {
3137         /*
3138         * We drop p_lock before grabbing the address
3139         * space lock in order to avoid a deadlock with
3140         * the clock thread. The process will not
3141         * disappear and its address space will not
3142         * change because it is marked P_PR_LOCK.
3143         */
3144         mutex_exit(&p->p_lock);
3145         AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
3146         seg = as_segat(as, va);
3147         if (seg != NULL &&
3148             seg->s_ops == &segvn_ops &&
3149             segop_getvp(seg, va, &xvp) == 0 &&
3149             SEGOP_GETVP(seg, va, &xvp) == 0 &&
3150             xvp != NULL &&
3151             xvp->v_type == VREG) {
3152             VN_HOLD(xvp);
3153         } else {
3154             error = EINVAL;
3155         }
3156         AS_LOCK_EXIT(as, &as->a_lock);
3157         mutex_enter(&p->p_lock);
3158     } else if ((xvp = p->p_exec) == NULL) {
3159         error = EINVAL;
3160     } else {
3161         VN_HOLD(xvp);
3162     }

3164     prunlock(pnp);

3166     if (error == 0) {
3167         if ((error = VOP_ACCESS(xvp, VREAD, 0, cr, NULL)) == 0)
3168             error = fassign(&xvp, FREAD, &n);
3169         if (error) {
3170             VN_RELE(xvp);
3171         } else {
3172             *rvalp = n;
3173         }
3174     }

```

```

3176         return (error);
3177     }
_____unchanged_portion_omitted_____

3511 /*
3512  * Return an array of structures with memory map information.
3513  * We allocate here; the caller must deallocate.
3514  * The caller is also responsible to append the zero-filled entry
3515  * that terminates the PIOCMAPOUT output buffer.
3516  */
3517 static int
3518 oprgetmap(proc_t *p, list_t *iolhead)
3519 {
3520     struct as *as = p->p_as;
3521     prmap_t *mp;
3522     struct seg *seg;
3523     struct seg *brkseg, *stkseg;
3524     uint_t prot;

3526     ASSERT(as != &kas && AS_WRITE_HELD(as, &as->a_lock));

3528     /*
3529     * Request an initial buffer size that doesn't waste memory
3530     * if the address space has only a small number of segments.
3531     */
3532     pr_iol_initlist(iolhead, sizeof (*mp), avl_numnodes(&as->a_segtree));

3534     if ((seg = AS_SEGFIRST(as)) == NULL)
3535         return (0);

3537     brkseg = break_seg(p);
3538     stkseg = as_segat(as, prgetstackbase(p));

3540     do {
3541         caddr_t eaddr = seg->s_base + pr_getsegsz(seg, 0);
3542         caddr_t saddr, naddr;
3543         void *tmp = NULL;

3545         for (saddr = seg->s_base; saddr < eaddr; saddr = naddr) {
3546             prot = pr_getprot(seg, 0, &tmp, &saddr, &naddr, eaddr);
3547             if (saddr == naddr)
3548                 continue;

3550             mp = pr_iol_newbuf(iolhead, sizeof (*mp));

3552             mp->pr_vaddr = saddr;
3553             mp->pr_size = naddr - saddr;
3554             mp->pr_off = segop_getoffset(seg, saddr);
3554             mp->pr_off = SEGOP_GETOFFSET(seg, saddr);
3555             mp->pr_mflags = 0;
3556             if (prot & PROT_READ)
3557                 mp->pr_mflags |= MA_READ;
3558             if (prot & PROT_WRITE)
3559                 mp->pr_mflags |= MA_WRITE;
3560             if (prot & PROT_EXEC)
3561                 mp->pr_mflags |= MA_EXEC;
3562             if (segop_gettype(seg, saddr) & MAP_SHARED)
3562             if (SEGOP_GETTYPE(seg, saddr) & MAP_SHARED)
3563                 mp->pr_mflags |= MA_SHARED;
3564             if (seg == brkseg)
3565                 mp->pr_mflags |= MA_BREAK;
3566             else if (seg == stkseg)
3567                 mp->pr_mflags |= MA_STACK;
3568             mp->pr_pagesize = PAGESIZE;
3569         }

```

```

3570         ASSERT(tmp == NULL);
3571     } while ((seg = AS_SEGNEXT(as, seg)) != NULL);

3573     return (0);
3574 }

3576 #ifdef _SYSCALL32_IMPL
3577 static int
3578 oprgetmap32(proc_t *p, list_t *iolhead)
3579 {
3580     struct as *as = p->p_as;
3581     ioc_prmap32_t *mp;
3582     struct seg *seg;
3583     struct seg *brkseg, *stkseg;
3584     uint_t prot;

3586     ASSERT(as != &kas && AS_WRITE_HELD(as, &as->a_lock));

3588     /*
3589      * Request an initial buffer size that doesn't waste memory
3590      * if the address space has only a small number of segments.
3591      */
3592     pr_iol_initlist(iolhead, sizeof (*mp), avl_numnodes(&as->a_segtree));

3594     if ((seg = AS_SEGFIRST(as)) == NULL)
3595         return (0);

3597     brkseg = break_seg(p);
3598     stkseg = as_segat(as, prgetstackbase(p));

3600     do {
3601         caddr_t eaddr = seg->s_base + pr_getsegsz(seg, 0);
3602         caddr_t saddr, naddr;
3603         void *tmp = NULL;

3605         for (saddr = seg->s_base; saddr < eaddr; saddr = naddr) {
3606             prot = pr_getprot(seg, 0, &tmp, &saddr, &naddr, eaddr);
3607             if (saddr == naddr)
3608                 continue;

3610             mp = pr_iol_newbuf(iolhead, sizeof (*mp));

3612             mp->pr_vaddr = (caddr32_t)(uintptr_t)saddr;
3613             mp->pr_size = (size32_t)(naddr - saddr);
3614             mp->pr_off = (off32_t)segop_getoffset(seg, saddr);
3615             mp->pr_off = (off32_t)SEGOP_GETOFFSET(seg, saddr);
3616             mp->pr_mflags = 0;
3617             if (prot & PROT_READ)
3618                 mp->pr_mflags |= MA_READ;
3619             if (prot & PROT_WRITE)
3620                 mp->pr_mflags |= MA_WRITE;
3621             if (prot & PROT_EXEC)
3622                 mp->pr_mflags |= MA_EXEC;
3623             if (segop_gettype(seg, saddr) & MAP_SHARED)
3624                 if (SEGOP_GETTYPE(seg, saddr) & MAP_SHARED)
3625                     mp->pr_mflags |= MA_SHARED;
3626             if (seg == brkseg)
3627                 mp->pr_mflags |= MA_BREAK;
3628             else if (seg == stkseg)
3629                 mp->pr_mflags |= MA_STACK;
3630             mp->pr_pagesize = PAGESIZE;
3631         }
3632     } while ((seg = AS_SEGNEXT(as, seg)) != NULL);

3633     return (0);

```

```

3634 }
3635     unchanged_portion_omitted
3698 #endif /* _SYSCALL32_IMPL */

3700 /*
3701  * Read old /proc page data information.
3702  */
3703 int
3704 oprpdread(struct as *as, uint_t hatid, struct uio *uiop)
3705 {
3706     caddr_t buf;
3707     size_t size;
3708     prpageheader_t *php;
3709     prasmap_t *pmp;
3710     struct seg *seg;
3711     int error;

3713     again:
3714     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);

3716     if ((seg = AS_SEGFIRST(as)) == NULL) {
3717         AS_LOCK_EXIT(as, &as->a_lock);
3718         return (0);
3719     }
3720     size = oprpdsize(as);
3721     if (uiop->uio_resid < size) {
3722         AS_LOCK_EXIT(as, &as->a_lock);
3723         return (E2BIG);
3724     }

3726     buf = kmem_zalloc(size, KM_SLEEP);
3727     php = (prpageheader_t *)buf;
3728     pmp = (prasmap_t *) (buf + sizeof (prpageheader_t));

3730     hrt2ts(gethrtime(), &php->pr_tstamp);
3731     php->pr_nmap = 0;
3732     php->pr_npage = 0;
3733     do {
3734         caddr_t eaddr = seg->s_base + pr_getsegsz(seg, 0);
3735         caddr_t saddr, naddr;
3736         void *tmp = NULL;

3738         for (saddr = seg->s_base; saddr < eaddr; saddr = naddr) {
3739             size_t len;
3740             size_t npage;
3741             uint_t prot;
3742             uintptr_t next;

3744             prot = pr_getprot(seg, 0, &tmp, &saddr, &naddr, eaddr);
3745             if ((len = naddr - saddr) == 0)
3746                 continue;
3747             npage = len / PAGESIZE;
3748             next = (uintptr_t)(pmp + 1) + roundup(npage);
3749             /*
3750              * It's possible that the address space can change
3751              * subtly even though we're holding as->a_lock
3752              * due to the nondeterminism of page_exists() in
3753              * the presence of asynchronously flushed pages or
3754              * mapped files whose sizes are changing.
3755              * page_exists() may be called indirectly from
3756              * pr_getprot() by a segop_incore() routine.
3757              * pr_getprot() by a SEGOP_INCORE() routine.
3758              * If this happens we need to make sure we don't
3759              * overrun the buffer whose size we computed based
3760              * on the initial iteration through the segments.
3761              * Once we've detected an overflow, we need to clean

```

```

3761     * up the temporary memory allocated in pr_getprot()
3762     * and retry. If there's a pending signal, we return
3763     * EINTR so that this thread can be dislodged if
3764     * a latent bug causes us to spin indefinitely.
3765     */
3766     if (next > (uintptr_t)buf + size) {
3767         pr_getprot_done(&tmp);
3768         AS_LOCK_EXIT(as, &as->a_lock);
3770
3771         kmem_free(buf, size);
3772
3773         if (ISSIG(curthread, JUSTLOOKING))
3774             return (EINTR);
3775
3776         goto again;
3777     }
3778     php->pr_nmap++;
3779     php->pr_npage += npage;
3780     pmp->pr_vaddr = saddr;
3781     pmp->pr_npage = npage;
3782     pmp->pr_off = segop_getoffset(seg, saddr);
3783     pmp->pr_off = SEGOP_GETOFFSET(seg, saddr);
3784     pmp->pr_mflags = 0;
3785     if (prot & PROT_READ)
3786         pmp->pr_mflags |= MA_READ;
3787     if (prot & PROT_WRITE)
3788         pmp->pr_mflags |= MA_WRITE;
3789     if (prot & PROT_EXEC)
3790         pmp->pr_mflags |= MA_EXEC;
3791     if (segop_gettype(seg, saddr) & MAP_SHARED)
3792         if (SEGOP_GETTYPE(seg, saddr) & MAP_SHARED)
3793             pmp->pr_mflags |= MA_SHARED;
3794     pmp->pr_pagesize = PAGE_SIZE;
3795     hat_getstat(as, saddr, len, hatid,
3796         (char *) (pmp + 1), HAT_SYNC_ZERORM);
3797     pmp = (prasm_t *) next;
3798     }
3799     ASSERT(tmp == NULL);
3800     } while ((seg = AS_SEGNEXT(as, seg)) != NULL);
3801
3802     AS_LOCK_EXIT(as, &as->a_lock);
3803
3804     ASSERT((uintptr_t)pmp <= (uintptr_t)buf + size);
3805     error = uiomove(buf, (caddr_t)pmp - buf, UIO_READ, uiop);
3806     kmem_free(buf, size);
3807
3808     return (error);
3809 }
3810
3811 #ifdef _SYS_CALL32_IMPL
3812 int
3813 oprpdread32(struct as *as, uint_t hatid, struct uio *uiop)
3814 {
3815     caddr_t buf;
3816     size_t size;
3817     ioc_prpageheader32_t *php;
3818     ioc_prasm32_t *pmp;
3819     struct seg *seg;
3820     int error;
3821
3822     again:
3823     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
3824
3825     if ((seg = AS_SEGFIRST(as)) == NULL) {
3826         AS_LOCK_EXIT(as, &as->a_lock);

```

```

3825         return (0);
3826     }
3827     size = oprpdsz32(as);
3828     if (uiop->uio_resid < size) {
3829         AS_LOCK_EXIT(as, &as->a_lock);
3830         return (E2BIG);
3831     }
3832
3833     buf = kmem_zalloc(size, KM_SLEEP);
3834     php = (ioc_prpageheader32_t *)buf;
3835     pmp = (ioc_prasm32_t *) (buf + sizeof (ioc_prpageheader32_t));
3836
3837     hrt2ts32(gethrtime(), &php->pr_tstamp);
3838     php->pr_nmap = 0;
3839     php->pr_npage = 0;
3840     do {
3841         caddr_t eaddr = seg->s_base + pr_getsegsz(seg, 0);
3842         caddr_t saddr, naddr;
3843         void *tmp = NULL;
3844
3845         for (saddr = seg->s_base; saddr < eaddr; saddr = naddr) {
3846             size_t len;
3847             size_t npage;
3848             uint_t prot;
3849             uintptr_t next;
3850
3851             prot = pr_getprot(seg, 0, &tmp, &saddr, &naddr, eaddr);
3852             if ((len = naddr - saddr) == 0)
3853                 continue;
3854             npage = len / PAGE_SIZE;
3855             next = (uintptr_t) (pmp + 1) + round4(npage);
3856             /*
3857              * It's possible that the address space can change
3858              * subtly even though we're holding as->a_lock
3859              * due to the nondeterminism of page_exists() in
3860              * the presence of asynchronously flushed pages or
3861              * mapped files whose sizes are changing.
3862              * page_exists() may be called indirectly from
3863              * pr_getprot() by a segop_incore() routine.
3864              * pr_getprot() by a segop_incore() routine.
3865              * If this happens we need to make sure we don't
3866              * overrun the buffer whose size we computed based
3867              * on the initial iteration through the segments.
3868              * Once we've detected an overflow, we need to clean
3869              * up the temporary memory allocated in pr_getprot()
3870              * and retry. If there's a pending signal, we return
3871              * EINTR so that this thread can be dislodged if
3872              * a latent bug causes us to spin indefinitely.
3873              */
3874             if (next > (uintptr_t)buf + size) {
3875                 pr_getprot_done(&tmp);
3876                 AS_LOCK_EXIT(as, &as->a_lock);
3877
3878                 kmem_free(buf, size);
3879
3880                 if (ISSIG(curthread, JUSTLOOKING))
3881                     return (EINTR);
3882
3883                 goto again;
3884             }
3885
3886             php->pr_nmap++;
3887             php->pr_npage += npage;
3888             pmp->pr_vaddr = (uint32_t) (uintptr_t) saddr;
3889             pmp->pr_npage = (uint32_t) npage;
3890             pmp->pr_off = (int32_t) segop_getoffset(seg, saddr);

```

```
3889     pmp->pr_off = (int32_t)SEGOP_GETOFFSET(seg, saddr);
3890     pmp->pr_mflags = 0;
3891     if (prot & PROT_READ)
3892         pmp->pr_mflags |= MA_READ;
3893     if (prot & PROT_WRITE)
3894         pmp->pr_mflags |= MA_WRITE;
3895     if (prot & PROT_EXEC)
3896         pmp->pr_mflags |= MA_EXEC;
3897     if (segop_gettype(seg, saddr) & MAP_SHARED)
3898         if (SEGOP_GETTYPE(seg, saddr) & MAP_SHARED)
3899             pmp->pr_mflags |= MA_SHARED;
3899     pmp->pr_pagesize = PAGE_SIZE;
3900     hat_getstat(as, saddr, len, hatid,
3901                (char *) (pmp + 1), HAT_SYNC_ZERORM);
3902     pmp = (ioc_prasmap32_t *) next;
3903 }
3904 ASSERT(tmp == NULL);
3905 } while ((seg = AS_SEGNEXT(as, seg)) != NULL);

3907 AS_LOCK_EXIT(as, &as->a_lock);

3909 ASSERT((uintptr_t)pmp == (uintptr_t)buf + size);
3910 error = uiomove(buf, (caddr_t)pmp - buf, UIO_READ, uiop);
3911 kmem_free(buf, size);

3913     return (error);
3914 }
```

_____unchanged_portion_omitted_____

```

*****
112686 Tue Nov 24 09:34:45 2015
new/usr/src/uts/common/fs/proc/prsubr.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

1475 struct seg *
1476 break_seg(proc_t *p)
1477 {
1478     caddr_t addr = p->p_brkbase;
1479     struct seg *seg;
1480     struct vnode *vp;

1482     if (p->p_brksize != 0)
1483         addr += p->p_brksize - 1;
1484     seg = as_segat(p->p_as, addr);
1485     if (seg != NULL && seg->s_ops == &segvn_ops &&
1486         (segop_getvp(seg, seg->s_base, &vp) != 0 || vp == NULL))
1487         (SEGOP_GETVP(seg, seg->s_base, &vp) != 0 || vp == NULL))
1488         return (seg);
1489     return (NULL);
_____unchanged_portion_omitted_____

1607 /*
1608  * Return an array of structures with memory map information.
1609  * We allocate here; the caller must deallocate.
1610  */
1611 int
1612 prgetmap(proc_t *p, int reserved, list_t *iolhead)
1613 {
1614     struct as *as = p->p_as;
1615     prmap_t *mp;
1616     struct seg *seg;
1617     struct seg *brkseg, *stkseg;
1618     struct vnode *vp;
1619     struct vattr vattr;
1620     uint_t prot;

1622     ASSERT(as != &kas && AS_WRITE_HELD(as, &as->a_lock));

1624     /*
1625      * Request an initial buffer size that doesn't waste memory
1626      * if the address space has only a small number of segments.
1627      */
1628     pr_iol_initlist(iolhead, sizeof (*mp), avl_numnodes(&as->a_segtree));

1630     if ((seg = AS_SEGFIRST(as)) == NULL)
1631         return (0);

1633     brkseg = break_seg(p);
1634     stkseg = as_segat(as, prgetstackbase(p));

1636     do {
1637         caddr_t eaddr = seg->s_base + pr_getsegsz(seg, reserved);
1638         caddr_t saddr, naddr;
1639         void *tmp = NULL;

1641         for (saddr = seg->s_base; saddr < eaddr; saddr = naddr) {
1642             prot = pr_getprot(seg, reserved, &tmp,
1643                 &saddr, &naddr, eaddr);
1644             if (saddr == naddr)
1645                 continue;

1647             mp = pr_iol_newbuf(iolhead, sizeof (*mp));

```

```

1649     mp->pr_vaddr = (uintptr_t)saddr;
1650     mp->pr_size = naddr - saddr;
1651     mp->pr_offset = segop_getoffset(seg, saddr);
1652     mp->pr_offset = SEGOP_GETOFFSET(seg, saddr);
1653     mp->pr_mflags = 0;
1654     if (prot & PROT_READ)
1655         mp->pr_mflags |= MA_READ;
1656     if (prot & PROT_WRITE)
1657         mp->pr_mflags |= MA_WRITE;
1658     if (prot & PROT_EXEC)
1659         mp->pr_mflags |= MA_EXEC;
1660     if (segop_gettype(seg, saddr) & MAP_SHARED)
1661         mp->pr_mflags |= MA_SHARED;
1662     if (segop_gettype(seg, saddr) & MAP_NORESERVE)
1663         mp->pr_mflags |= MA_NORESERVE;
1664     if (seg->s_ops == &segspt_shmops ||
1665         (seg->s_ops == &segvn_ops &&
1666             (segop_getvp(seg, saddr, &vp) != 0 || vp == NULL)))
1667         (SEGOP_GETVP(seg, saddr, &vp) != 0 || vp == NULL))
1668         mp->pr_mflags |= MA_ANON;
1669     if (seg == brkseg)
1670         mp->pr_mflags |= MA_BREAK;
1671     else if (seg == stkseg) {
1672         mp->pr_mflags |= MA_STACK;
1673         if (reserved) {
1674             size_t maxstack =
1675                 ((size_t)p->p_stk_ctl +
1676                 PAGEOFFSET) & PAGEMASK;
1677             mp->pr_vaddr =
1678                 (uintptr_t)prgetstackbase(p) +
1679                 p->p_stksize - maxstack;
1680             mp->pr_size = (uintptr_t)naddr -
1681                 mp->pr_vaddr;
1682         }
1683     }
1684     if (seg->s_ops == &segspt_shmops)
1685         mp->pr_mflags |= MA_ISM | MA_SHM;
1686     mp->pr_pagesize = PAGESIZE;

1688     /*
1689      * Manufacture a filename for the "object" directory.
1690      */
1691     vattr.va_mask = AT_FSID|AT_NODEID;
1692     if (seg->s_ops == &segvn_ops &&
1693         (segop_getvp(seg, saddr, &vp) == 0 &&
1694             SEGOP_GETVP(seg, saddr, &vp) == 0 &&
1695             vp != NULL && vp->v_type == VREG &&
1696             VOP_GETATTR(vp, &vattr, 0, CRED(), NULL) == 0) {
1697         if (vp == p->p_exec)
1698             (void) strcpy(mp->pr_mapname, "a.out");
1699         else
1700             pr_object_name(mp->pr_mapname,
1701                 vp, &vattr);
1702     }

1704     /*
1705      * Get the SysV shared memory id, if any.
1706      */
1707     if ((mp->pr_mflags & MA_SHARED) && p->p_segacct &&
1708         (mp->pr_shmid = shmgetid(p, seg->s_base)) !=
1709         SHMID_NONE) {
1710         if (mp->pr_shmid == SHMID_FREE)
1711             mp->pr_shmid = -1;

```

```

1710         mp->pr_mflags |= MA_SHM;
1711     } else {
1712         mp->pr_shmid = -1;
1713     }
1714 }
1715     ASSERT(tmp == NULL);
1716 } while ((seg = AS_SEGNEXT(as, seg)) != NULL);

1718     return (0);
1719 }

1721 #ifdef _SYSCALL32_IMPL
1722 int
1723 prgetmap32(proc_t *p, int reserved, list_t *iolhead)
1724 {
1725     struct as *as = p->p_as;
1726     prmap32_t *mp;
1727     struct seg *seg;
1728     struct seg *brkseg, *stkseg;
1729     struct vnode *vp;
1730     struct vattr vattr;
1731     uint_t prot;

1733     ASSERT(as != &kas && AS_WRITE_HELD(as, &as->a_lock));

1735     /*
1736      * Request an initial buffer size that doesn't waste memory
1737      * if the address space has only a small number of segments.
1738      */
1739     pr_iol_initlist(iolhead, sizeof (*mp), avl_numnodes(&as->a_segtree));

1741     if ((seg = AS_SEGFIRST(as)) == NULL)
1742         return (0);

1744     brkseg = break_seg(p);
1745     stkseg = as_segat(as, prgetstackbase(p));

1747     do {
1748         caddr_t eaddr = seg->s_base + pr_getsegsz(seg, reserved);
1749         caddr_t saddr, naddr;
1750         void *tmp = NULL;

1752         for (saddr = seg->s_base; saddr < eaddr; saddr = naddr) {
1753             prot = pr_getprot(seg, reserved, &tmp,
1754                 &saddr, &naddr, eaddr);
1755             if (saddr == naddr)
1756                 continue;

1758             mp = pr_iol_newbuf(iolhead, sizeof (*mp));

1760             mp->pr_vaddr = (caddr32_t)(uintptr_t)saddr;
1761             mp->pr_size = (size32_t)(naddr - saddr);
1762             mp->pr_offset = segop_getoffset(seg, saddr);
1763             mp->pr_offset = SEGOP_GETOFFSET(seg, saddr);
1764             mp->pr_mflags = 0;
1765             if (prot & PROT_READ)
1766                 mp->pr_mflags |= MA_READ;
1767             if (prot & PROT_WRITE)
1768                 mp->pr_mflags |= MA_WRITE;
1769             if (prot & PROT_EXEC)
1770                 mp->pr_mflags |= MA_EXEC;
1771             if (segop_gettype(seg, saddr) & MAP_SHARED)
1772                 mp->pr_mflags |= MA_SHARED;
1773             if (SEGOP_GETTYPE(seg, saddr) & MAP_SHARED)
1774                 mp->pr_mflags |= MA_SHARED;
1775             if (segop_gettype(seg, saddr) & MAP_NORESERVE)

```

```

1776         if (SEGOP_GETTYPE(seg, saddr) & MAP_NORESERVE)
1777             mp->pr_mflags |= MA_NORESERVE;
1778         if (seg->s_ops == &segspt_shmops ||
1779             (seg->s_ops == &segvn_ops &&
1780             (segop_getvp(seg, saddr, &vp) != 0 || vp == NULL)))
1781             (SEGOP_GETVP(seg, saddr, &vp) != 0 || vp == NULL))
1782             mp->pr_mflags |= MA_ANON;
1783         if (seg == brkseg)
1784             mp->pr_mflags |= MA_BREAK;
1785         else if (seg == stkseg) {
1786             mp->pr_mflags |= MA_STACK;
1787             if (reserved) {
1788                 size_t maxstack =
1789                     ((size_t)p->p_stk_ctl +
1790                     PAGEOFFSET) & PAGEMASK;
1791                 uintptr_t vaddr =
1792                     (uintptr_t)prgetstackbase(p) +
1793                     p->p_stksize - maxstack;
1794                 mp->pr_vaddr = (caddr32_t)vaddr;
1795                 mp->pr_size = (size32_t)
1796                     ((uintptr_t)naddr - vaddr);
1797             }
1798         }
1799         if (seg->s_ops == &segspt_shmops)
1800             mp->pr_mflags |= MA_ISM | MA_SHM;
1801         mp->pr_pagesize = PAGE_SIZE;

1803         /*
1804          * Manufacture a filename for the "object" directory.
1805          */
1806         vattr.va_mask = AT_FSID|AT_NODEID;
1807         if (seg->s_ops == &segvn_ops &&
1808             segop_getvp(seg, saddr, &vp) == 0 &&
1809             SEGOP_GETVP(seg, saddr, &vp) == 0 &&
1810             vp != NULL && vp->v_type == VREG &&
1811             VOP_GETATTR(vp, &vattr, 0, CRED(), NULL) == 0) {
1812             if (vp == p->p_exec)
1813                 (void) strcpy(mp->pr_mapname, "a.out");
1814             else
1815                 pr_object_name(mp->pr_mapname,
1816                     vp, &vattr);
1817         }

1819         /*
1820          * Get the SysV shared memory id, if any.
1821          */
1822         if ((mp->pr_mflags & MA_SHARED) && p->p_segacct &&
1823             (mp->pr_shmid = shmgetid(p, seg->s_base)) !=
1824             SHMID_NONE) {
1825             if (mp->pr_shmid == SHMID_FREE)
1826                 mp->pr_shmid = -1;

1828             mp->pr_mflags |= MA_SHM;
1829         } else {
1830             mp->pr_shmid = -1;
1831         }
1832     }
1833     ASSERT(tmp == NULL);
1834 } while ((seg = AS_SEGNEXT(as, seg)) != NULL);

1836     return (0);
1837 }

1839 #endif /* _SYSCALL32_IMPL */
1840 */

```



```

1898 * Read page data information.
1899 */
1900 int
1901 prp_dread(proc_t *p, uint_t hatid, struct uio *uiop)
1902 {
1903     struct as *as = p->p_as;
1904     caddr_t buf;
1905     size_t size;
1906     prpageheader_t *php;
1907     prasmap_t *pmp;
1908     struct seg *seg;
1909     int error;

1911 again:
1912     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);

1914     if ((seg = AS_SEGFIRST(as)) == NULL) {
1915         AS_LOCK_EXIT(as, &as->a_lock);
1916         return (0);
1917     }
1918     size = prpdsiz(eas);
1919     if (uiop->uio_resid < size) {
1920         AS_LOCK_EXIT(as, &as->a_lock);
1921         return (E2BIG);
1922     }

1924     buf = kmem_zalloc(size, KM_SLEEP);
1925     php = (prpageheader_t *)buf;
1926     pmp = (prasmap_t *) (buf + sizeof (prpageheader_t));

1928     hrt2ts(gethrtime(), &php->pr_tstamp);
1929     php->pr_nmap = 0;
1930     php->pr_npage = 0;
1931     do {
1932         caddr_t eaddr = seg->s_base + pr_getsegsiz(e, seg, 0);
1933         caddr_t saddr, naddr;
1934         void *tmp = NULL;

1936         for (saddr = seg->s_base; saddr < eaddr; saddr = naddr) {
1937             struct vnode *vp;
1938             struct vattr vattr;
1939             size_t len;
1940             size_t npage;
1941             uint_t prot;
1942             uintptr_t next;

1944             prot = pr_getprot(seg, 0, &tmp, &saddr, &naddr, eaddr);
1945             if ((len = (size_t)(naddr - saddr)) == 0)
1946                 continue;
1947             npage = len / PAGE_SIZE;
1948             next = (uintptr_t)(pmp + 1) + round8(npage);
1949             /*
1950              * It's possible that the address space can change
1951              * subtly even though we're holding as->a_lock
1952              * due to the nondeterminism of page_exists() in
1953              * the presence of asynchronously flushed pages or
1954              * mapped files whose sizes are changing.
1955              * page_exists() may be called indirectly from
1956              * pr_getprot() by a segop_incore() routine.
1957              * pr_getprot() by a SEGOP_INCORE() routine.
1958              * If this happens we need to make sure we don't
1959              * overrun the buffer whose size we computed based
1960              * on the initial iteration through the segments.
1961              * Once we've detected an overflow, we need to clean
1962              * up the temporary memory allocated in pr_getprot()
1963              * and retry. If there's a pending signal, we return

```

```

1963         * EINTR so that this thread can be dislodged if
1964         * a latent bug causes us to spin indefinitely.
1965         */
1966         if (next > (uintptr_t)buf + size) {
1967             pr_getprot_done(&tmp);
1968             AS_LOCK_EXIT(as, &as->a_lock);

1970             kmem_free(buf, size);

1972             if (ISSIG(curthread, JUSTLOOKING))
1973                 return (EINTR);

1975             goto again;
1976         }

1978         php->pr_nmap++;
1979         php->pr_npage += npage;
1980         pmp->pr_vaddr = (uintptr_t)saddr;
1981         pmp->pr_npage = npage;
1982         pmp->pr_offset = segop_getoffset(seg, saddr);
1983         pmp->pr_offset = SEGOP_GETOFFSET(seg, saddr);
1984         pmp->pr_mflags = 0;
1985         if (prot & PROT_READ)
1986             pmp->pr_mflags |= MA_READ;
1987         if (prot & PROT_WRITE)
1988             pmp->pr_mflags |= MA_WRITE;
1989         if (prot & PROT_EXEC)
1990             pmp->pr_mflags |= MA_EXEC;
1991         if (segop_gettype(seg, saddr) & MAP_SHARED)
1992             pmp->pr_mflags |= MA_SHARED;
1993         if (segop_gettype(seg, saddr) & MAP_NORESERVE)
1994             pmp->pr_mflags |= MA_NORESERVE;
1995         if (seg->s_ops == &segspt_shmops ||
1996             (seg->s_ops == &segvn_ops &&
1997              (segop_getvp(seg, saddr, &vp) != 0 || vp == NULL)))
1998             pmp->pr_mflags |= MA_ANON;
1999         if (seg->s_ops == &segspt_shmops)
2000             pmp->pr_mflags |= MA_ISM | MA_SHM;
2001         pmp->pr_pagesize = PAGE_SIZE;
2002         /*
2003          * Manufacture a filename for the "object" directory.
2004          */
2005         vattr.va_mask = AT_FSID|AT_NODEID;
2006         if (seg->s_ops == &segvn_ops &&
2007             segop_getvp(seg, saddr, &vp) == 0 &&
2008             SEGOP_GETVP(seg, saddr, &vp) == 0 &&
2009             vp != NULL && vp->v_type == VREG &&
2010             VOP_GETATTR(vp, &vattr, 0, CRED(), NULL) == 0) {
2011             if (vp == p->p_exec)
2012                 (void) strcpy(pmp->pr_mapname, "a.out");
2013             else
2014                 pr_object_name(pmp->pr_mapname,
2015                               vp, &vattr);
2016         }
2017         /*
2018          * Get the SysV shared memory id, if any.
2019          */
2020         if ((pmp->pr_mflags & MA_SHARED) && p->p_segacct &&
2021             (pmp->pr_shmid = shmgetid(p, seg->s_base)) !=
2022             SHMID_NONE) {
2023             if (pmp->pr_shmid == SHMID_FREE)
2024                 pmp->pr_shmid = -1;

```

```

2025         pmp->pr_mflags |= MA_SHM;
2026     } else {
2027         pmp->pr_shmid = -1;
2028     }
2030     hat_getstat(as, saddr, len, hatid,
2031         (char *) (pmp + 1), HAT_SYNC_ZERORM);
2032     pmp = (prasmap_t *) next;
2033 }
2034     ASSERT(tmp == NULL);
2035 } while ((seg = AS_SEGNEXT(as, seg)) != NULL);
2037 AS_LOCK_EXIT(as, &as->a_lock);
2039 ASSERT((uintptr_t)pmp <= (uintptr_t)buf + size);
2040 error = uiomove(buf, (caddr_t)pmp - buf, UIO_READ, uiop);
2041 kmem_free(buf, size);
2043     return (error);
2044 }
2046 #ifdef _SYSCALL32_IMPL
2047 int
2048 prpdread32(proc_t *p, uint_t hatid, struct uio *uiop)
2049 {
2050     struct as *as = p->p_as;
2051     caddr_t buf;
2052     size_t size;
2053     prpageheader32_t *php;
2054     prasmap32_t *pmp;
2055     struct seg *seg;
2056     int error;
2058 again:
2059     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
2061     if ((seg = AS_SEGFIRST(as)) == NULL) {
2062         AS_LOCK_EXIT(as, &as->a_lock);
2063         return (0);
2064     }
2065     size = prpdsize32(as);
2066     if (uiop->uio_resid < size) {
2067         AS_LOCK_EXIT(as, &as->a_lock);
2068         return (E2BIG);
2069     }
2071     buf = kmem_zalloc(size, KM_SLEEP);
2072     php = (prpageheader32_t *) buf;
2073     pmp = (prasmap32_t *) (buf + sizeof (prpageheader32_t));
2075     hrt2ts32(gethrtime(), &php->pr_tstamp);
2076     php->pr_nmap = 0;
2077     php->pr_npage = 0;
2078     do {
2079         caddr_t eaddr = seg->s_base + pr_getsegsz(seg, 0);
2080         caddr_t saddr, naddr;
2081         void *tmp = NULL;
2083         for (saddr = seg->s_base; saddr < eaddr; saddr = naddr) {
2084             struct vnode *vp;
2085             struct vattr vattr;
2086             size_t len;
2087             size_t npage;
2088             uint_t prot;
2089             uintptr_t next;

```

```

2091     prot = pr_getprot(seg, 0, &tmp, &saddr, &naddr, eaddr);
2092     if ((len = (size_t)(naddr - saddr)) == 0)
2093         continue;
2094     npage = len / PAGE_SIZE;
2095     next = (uintptr_t)(pmp + 1) + round8(npage);
2096     /*
2097      * It's possible that the address space can change
2098      * subtly even though we're holding as->a_lock
2099      * due to the nondeterminism of page_exists() in
2100      * the presence of asynchronously flushed pages or
2101      * mapped files whose sizes are changing.
2102      * page_exists() may be called indirectly from
2103      * pr_getprot() by a segop_incore() routine.
2104      * pr_getprot() by a SEGOP_INCORE() routine.
2105      * If this happens we need to make sure we don't
2106      * overrun the buffer whose size we computed based
2107      * on the initial iteration through the segments.
2108      * Once we've detected an overflow, we need to clean
2109      * up the temporary memory allocated in pr_getprot()
2110      * and retry. If there's a pending signal, we return
2111      * EINTR so that this thread can be dislodged if
2112      * a latent bug causes us to spin indefinitely.
2113      */
2114     if (next > (uintptr_t)buf + size) {
2115         pr_getprot_done(&tmp);
2116         AS_LOCK_EXIT(as, &as->a_lock);
2117         kmem_free(buf, size);
2119         if (ISSIG(curthread, JUSTLOOKING))
2120             return (EINTR);
2122         goto again;
2123     }
2125     php->pr_nmap++;
2126     php->pr_npage += npage;
2127     pmp->pr_vaddr = (caddr32_t)(uintptr_t)saddr;
2128     pmp->pr_npage = (size32_t)npage;
2129     pmp->pr_offset = segop_getoffset(seg, saddr);
2130     pmp->pr_offset = SEGOP_GETOFFSET(seg, saddr);
2131     pmp->pr_mflags = 0;
2132     if (prot & PROT_READ)
2133         pmp->pr_mflags |= MA_READ;
2134     if (prot & PROT_WRITE)
2135         pmp->pr_mflags |= MA_WRITE;
2136     if (prot & PROT_EXEC)
2137         pmp->pr_mflags |= MA_EXEC;
2138     if (segop_gettype(seg, saddr) & MAP_SHARED)
2139         pmp->pr_mflags |= MA_SHARED;
2140     if (segop_gettype(seg, saddr) & MAP_NORESERVE)
2141         pmp->pr_mflags |= MA_NORESERVE;
2142     if (seg->s_ops == &segspt_shmops ||
2143         (seg->s_ops == &segvn_ops &&
2144         (segop_getvp(seg, saddr, &vp) != 0 || vp == NULL)))
2145         pmp->pr_mflags |= MA_ANON;
2146     if (seg->s_ops == &segspt_shmops)
2147         pmp->pr_mflags |= MA_ISM | MA_SHM;
2148     pmp->pr_pagesize = PAGE_SIZE;
2149     /*
2150      * Manufacture a filename for the "object" directory.
2151      */

```

```

2151     vattr.va_mask = AT_FSID|AT_NODEID;
2152     if (seg->s_ops == &segvn_ops &&
2153         segop_getvp(seg, saddr, &vp) == 0 &&
2154         SEGOP_GETVP(seg, saddr, &vp) == 0 &&
2155         vp != NULL && vp->v_type == VREG &&
2156         VOP_GETATTR(vp, &vattr, 0, CRED(), NULL) == 0) {
2157         if (vp == p->p_exec)
2158             (void) strcpy(pmp->pr_mapname, "a.out");
2159         else
2160             pr_object_name(pmp->pr_mapname,
2161                 vp, &vattr);
2162     }
2163     /*
2164     * Get the SysV shared memory id, if any.
2165     */
2166     if ((pmp->pr_mflags & MA_SHARED) && p->p_segacct &&
2167         (pmp->pr_shmid = shmgetid(p, seg->s_base)) !=
2168         SHMID_NONE) {
2169         if (pmp->pr_shmid == SHMID_FREE)
2170             pmp->pr_shmid = -1;
2171     } else {
2172         pmp->pr_mflags |= MA_SHM;
2173     }
2174     pmp->pr_shmid = -1;
2175 }
2176
2177     hat_getstat(as, saddr, len, hatid,
2178         (char *) (pmp + 1), HAT_SYNC_ZERORM);
2179     pmp = (prasm32_t *) next;
2180 }
2181     ASSERT(tmp == NULL);
2182 } while ((seg = AS_SEGNEXT(as, seg)) != NULL);
2183
2184     AS_LOCK_EXIT(as, &as->a_lock);
2185
2186     ASSERT((uintptr_t)pmp <= (uintptr_t)buf + size);
2187     error = uiomove(buf, (caddr_t)pmp - buf, UIO_READ, uiop);
2188     kmem_free(buf, size);
2189
2190     return (error);
2191 }

```

unchanged portion omitted

```

3299 /*
3300 * This one is called by the traced process to unwatch all the
3301 * pages while deallocating the list of watched_page structs.
3302 */
3303 void
3304 pr_free_watched_pages(proc_t *p)
3305 {
3306     struct as *as = p->p_as;
3307     struct watched_page *pwp;
3308     uint_t prot;
3309     int retrycnt, err;
3310     void *cookie;
3311
3312     if (as == NULL || avl_numnodes(&as->a_wpage) == 0)
3313         return;
3314
3315     ASSERT(MUTEX_NOT_HELD(&curproc->p_lock));
3316     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
3317
3318     pwp = avl_first(&as->a_wpage);
3319
3320     cookie = NULL;

```

```

3321     while ((pwp = avl_destroy_nodes(&as->a_wpage, &cookie)) != NULL) {
3322         retrycnt = 0;
3323         if ((prot = pwp->wp_oprot) != 0) {
3324             caddr_t addr = pwp->wp_vaddr;
3325             struct seg *seg;
3326
3327             retry:
3328                 if ((pwp->wp_prot != prot ||
3329                     (pwp->wp_flags & WP_NOWATCH)) &&
3330                     (seg = as_segat(as, addr)) != NULL) {
3331                     err = segop_setprot(seg, addr, PAGE_SIZE, prot);
3332                     err = SEGOP_SETPROT(seg, addr, PAGE_SIZE, prot);
3333                     if (err == IE_RETRY) {
3334                         ASSERT(retrycnt == 0);
3335                         retrycnt++;
3336                         goto retry;
3337                     }
3338                 }
3339             kmem_free(pwp, sizeof (struct watched_page));
3340         }
3341     }
3342     avl_destroy(&as->a_wpage);
3343     p->p_wprot = NULL;
3344
3345     AS_LOCK_EXIT(as, &as->a_lock);
3346 }
3347
3348 /*
3349 * Insert a watched area into the list of watched pages.
3350 * If oflags is zero then we are adding a new watched area.
3351 * Otherwise we are changing the flags of an existing watched area.
3352 */
3353 static int
3354 set_watched_page(proc_t *p, caddr_t vaddr, caddr_t eaddr,
3355     ulong_t flags, ulong_t oflags)
3356 {
3357     struct as *as = p->p_as;
3358     avl_tree_t *pwp_tree;
3359     struct watched_page *pwp, *newpwp;
3360     struct watched_page tpw;
3361     avl_index_t where;
3362     struct seg *seg;
3363     uint_t prot;
3364     caddr_t addr;
3365
3366     /*
3367     * We need to pre-allocate a list of structures before we grab the
3368     * address space lock to avoid calling kmem_alloc(KM_SLEEP) with locks
3369     * held.
3370     */
3371     newpwp = NULL;
3372     for (addr = (caddr_t)((uintptr_t)vaddr & (uintptr_t)PAGEMASK);
3373         addr < eaddr; addr += PAGE_SIZE) {
3374         pwp = kmem_zalloc(sizeof (struct watched_page), KM_SLEEP);
3375         pwp->wp_list = newpwp;
3376         newpwp = pwp;
3377     }
3378
3379     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
3380
3381     /*
3382     * Search for an existing watched page to contain the watched area.
3383     * If none is found, grab a new one from the available list
3384     * and insert it in the active list, keeping the list sorted
3385     * by user-level virtual address.

```

```

3386  */
3387  if (p->p_flag & SVFWAIT)
3388      pwp_tree = &p->p_wpage;
3389  else
3390      pwp_tree = &as->a_wpage;

3392 again:
3393  if (avl_numnodes(pwp_tree) > prnwatch) {
3394      AS_LOCK_EXIT(as, &as->a_lock);
3395      while (newpwp != NULL) {
3396          pwp = newpwp->wp_list;
3397          kmem_free(newpwp, sizeof (struct watched_page));
3398          newpwp = pwp;
3399      }
3400      return (E2BIG);
3401  }

3403  tpw.wp_vaddr = (caddr_t)((uintptr_t)vaddr & (uintptr_t)PAGEMASK);
3404  if ((pwp = avl_find(pwp_tree, &tpw, &where)) == NULL) {
3405      pwp = newpwp;
3406      newpwp = newpwp->wp_list;
3407      pwp->wp_list = NULL;
3408      pwp->wp_vaddr = (caddr_t)((uintptr_t)vaddr &
3409          (uintptr_t)PAGEMASK);
3410      avl_insert(pwp_tree, pwp, where);
3411  }

3413  ASSERT(vaddr >= pwp->wp_vaddr && vaddr < pwp->wp_vaddr + PAGE_SIZE);

3415  if (oflags & WA_READ)
3416      pwp->wp_read--;
3417  if (oflags & WA_WRITE)
3418      pwp->wp_write--;
3419  if (oflags & WA_EXEC)
3420      pwp->wp_exec--;

3422  ASSERT(pwp->wp_read >= 0);
3423  ASSERT(pwp->wp_write >= 0);
3424  ASSERT(pwp->wp_exec >= 0);

3426  if (flags & WA_READ)
3427      pwp->wp_read++;
3428  if (flags & WA_WRITE)
3429      pwp->wp_write++;
3430  if (flags & WA_EXEC)
3431      pwp->wp_exec++;

3433  if (!(p->p_flag & SVFWAIT)) {
3434      vaddr = pwp->wp_vaddr;
3435      if (pwp->wp_oprot == 0 &&
3436          (seg = as_segat(as, vaddr)) != NULL) {
3437          (void) segop_getprot(seg, vaddr, 0, &prot);
3437          SEGOP_GETPROT(seg, vaddr, 0, &prot);
3438          pwp->wp_oprot = (uchar_t)prot;
3439          pwp->wp_prot = (uchar_t)prot;
3440      }
3441      if (pwp->wp_oprot != 0) {
3442          prot = pwp->wp_oprot;
3443          if (pwp->wp_read)
3444              prot &= ~(PROT_READ|PROT_WRITE|PROT_EXEC);
3445          if (pwp->wp_write)
3446              prot &= ~PROT_WRITE;
3447          if (pwp->wp_exec)
3448              prot &= ~(PROT_READ|PROT_WRITE|PROT_EXEC);
3449          if (!(pwp->wp_flags & WP_NOWATCH) &&
3450              pwp->wp_prot != prot &&

```

```

3451          (pwp->wp_flags & WP_SETPROT) == 0) {
3452              pwp->wp_flags |= WP_SETPROT;
3453              pwp->wp_list = p->p_wprot;
3454              p->p_wprot = pwp;
3455          }
3456          pwp->wp_prot = (uchar_t)prot;
3457      }
3458  }

3460  /*
3461   * If the watched area extends into the next page then do
3462   * it over again with the virtual address of the next page.
3463   */
3464  if ((vaddr = pwp->wp_vaddr + PAGE_SIZE) < eaddr)
3465      goto again;

3467  AS_LOCK_EXIT(as, &as->a_lock);

3469  /*
3470   * Free any pages we may have over-allocated
3471   */
3472  while (newpwp != NULL) {
3473      pwp = newpwp->wp_list;
3474      kmem_free(newpwp, sizeof (struct watched_page));
3475      newpwp = pwp;
3476  }

3478  return (0);
3479 }

_____unchanged_portion_omitted_____

3614 static caddr_t
3615 pr_pagev_fill(prpagev_t *pagev, struct seg *seg, caddr_t addr, caddr_t eaddr)
3616 {
3617     ulong_t lastpg = seg_page(seg, eaddr - 1);
3618     ulong_t pn, pnlim;
3619     caddr_t saddr;
3620     size_t len;

3622     ASSERT(addr >= seg->s_base && addr <= eaddr);

3624     if (addr == eaddr)
3625         return (eaddr);

3627 refill:
3628     ASSERT(addr < eaddr);
3629     pagev->pg_pnbase = seg_page(seg, addr);
3630     pnlim = pagev->pg_pnbase + pagev->pg_npages;
3631     saddr = addr;

3633     if (lastpg < pnlim)
3634         len = (size_t)(eaddr - addr);
3635     else
3636         len = pagev->pg_npages * PAGE_SIZE;

3638     if (pagev->pg_incore != NULL) {
3639         /*
3640          * INCORE cleverly has different semantics than GETPROT:
3641          * it returns info on pages up to but NOT including addr + len.
3642          */
3643         (void) segop_incore(seg, addr, len, pagev->pg_incore);
3643         SEGOP_INCORE(seg, addr, len, pagev->pg_incore);
3644         pn = pagev->pg_pnbase;

3646         do {
3647             /*

```

```

3648         * Guilty knowledge here: We know that segvn_incore
3649         * returns more than just the low-order bit that
3650         * indicates the page is actually in memory. If any
3651         * bits are set, then the page has backing store.
3652         */
3653         if (pagev->pg_incore[pn++ - pagev->pg_pnbase])
3654             goto out;
3655
3656     } while ((addr += PAGE_SIZE) < eaddr && pn < pnlim);
3657
3658     /*
3659     * If we examined all the pages in the vector but we're not
3660     * at the end of the segment, take another lap.
3661     */
3662     if (addr < eaddr)
3663         goto refill;
3664 }
3665
3666 /*
3667 * Need to take len - 1 because addr + len is the address of the
3668 * first byte of the page just past the end of what we want.
3669 */
3670 out:
3671 (void) segop_getprot(seg, saddr, len - 1, pagev->pg_prot);
3672 SEGOP_GETPROT(seg, saddr, len - 1, pagev->pg_prot);
3673 return (addr);
3674 }
3675
3676 _____ unchanged_portion_omitted _____
3677
3678 size_t
3679 pr_getsegsz(struct seg *seg, int reserved)
3680 {
3681     size_t size = seg->s_size;
3682
3683     /*
3684     * If we're interested in the reserved space, return the size of the
3685     * segment itself. Everything else in this function is a special case
3686     * to determine the actual underlying size of various segment types.
3687     */
3688     if (reserved)
3689         return (size);
3690
3691     /*
3692     * If this is a segvn mapping of a regular file, return the smaller
3693     * of the segment size and the remaining size of the file beyond
3694     * the file offset corresponding to seg->s_base.
3695     */
3696     if (seg->s_ops == &segvn_ops) {
3697         vattr_t vattr;
3698         vnode_t *vp;
3699
3700         vattr.va_mask = AT_SIZE;
3701
3702         if (segop_getvp(seg, seg->s_base, &vp) == 0 &&
3703             if (SEGOP_GETVP(seg, seg->s_base, &vp) == 0 &&
3704                 vp != NULL && vp->v_type == VREG &&
3705                 VOP_GETATTR(vp, &vattr, 0, CRED(), NULL) == 0) {
3706             u_offset_t fsize = vattr.va_size;
3707             u_offset_t offset = segop_getoffset(seg, seg->s_base);
3708             u_offset_t offset = SEGOP_GETOFFSET(seg, seg->s_base);
3709
3710             if (fsize < offset)
3711                 fsize = 0;
3712             else
3713                 fsize -= offset;

```

```

3804         fsize = roundup(fsize, (u_offset_t)PAGE_SIZE);
3805
3806         if (fsize < (u_offset_t)size)
3807             size = (size_t)fsize;
3808     }
3809
3810     return (size);
3811 }
3812
3813 /*
3814 * If this is an ISM shared segment, don't include pages that are
3815 * beyond the real size of the spt segment that backs it.
3816 */
3817 if (seg->s_ops == &segspt_shmops)
3818     return (MIN(spt_realsize(seg), size));
3819
3820 /*
3821 * If this segment is a mapping from /dev/null, then this is a
3822 * reservation of virtual address space and has no actual size.
3823 * Such segments are backed by segdev and have type set to neither
3824 * MAP_SHARED nor MAP_PRIVATE.
3825 */
3826 if (seg->s_ops == &segdev_ops &&
3827     ((segop_gettype(seg, seg->s_base) &
3828      ((SEGOP_GETTYPE(seg, seg->s_base) &
3829       (MAP_SHARED | MAP_PRIVATE)) == 0))
3830     return (0);
3831
3832 /*
3833 * If this segment doesn't match one of the special types we handle,
3834 * just return the size of the segment itself.
3835 */
3836 return (size);
3837 }
3838
3839 _____ unchanged_portion_omitted _____
3840
3841 /*
3842 * Return an array of structures with extended memory map information.
3843 * We allocate here; the caller must deallocate.
3844 */
3845 int
3846 pr_getxmap(proc_t *p, list_t *iolhead)
3847 {
3848     struct as *as = p->p_as;
3849     prxmap_t *mp;
3850     struct seg *seg;
3851     struct seg *brkseg, *stkseg;
3852     struct vnode *vp;
3853     struct vattr vattr;
3854     uint_t prot;
3855
3856     ASSERT(as != &kas && AS_WRITE_HELD(as, &as->a_lock));
3857
3858     /*
3859     * Request an initial buffer size that doesn't waste memory
3860     * if the address space has only a small number of segments.
3861     */
3862     pr_iol_initlist(iolhead, sizeof (*mp), avl_numnodes(&as->a_segtree));
3863
3864     if ((seg = AS_SEGFIRST(as)) == NULL)
3865         return (0);
3866
3867     brkseg = break_seg(p);
3868     stkseg = as_segat(as, prgetstackbase(p));

```

```

4026     do {
4027         caddr_t eaddr = seg->s_base + pr_getsegsz(seg, 0);
4028         caddr_t saddr, naddr, baddr;
4029         void *tmp = NULL;
4030         ssize_t psz;
4031         char *parr;
4032         uint64_t npages;
4033         uint64_t pagenum;

4035     /*
4036      * Segment loop part one: iterate from the base of the segment
4037      * to its end, pausing at each address boundary (baddr) between
4038      * ranges that have different virtual memory protections.
4039      */
4040     for (saddr = seg->s_base; saddr < eaddr; saddr = baddr) {
4041         prot = pr_getprot(seg, 0, &tmp, &saddr, &baddr, eaddr);
4042         ASSERT(baddr >= saddr && baddr <= eaddr);

4044         /*
4045          * Segment loop part two: iterate from the current
4046          * position to the end of the protection boundary,
4047          * pausing at each address boundary (naddr) between
4048          * ranges that have different underlying page sizes.
4049          */
4050         for (; saddr < baddr; saddr = naddr) {
4051             psz = pr_getpagesize(seg, saddr, &naddr, baddr);
4052             ASSERT(naddr >= saddr && naddr <= baddr);

4054             mp = pr_iol_newbuf(iolhead, sizeof (*mp));

4056             mp->pr_vaddr = (uintptr_t)saddr;
4057             mp->pr_size = naddr - saddr;
4058             mp->pr_offset = segop_getoffset(seg, saddr);
4059             mp->pr_offset = SEGOP_GETOFFSET(seg, saddr);
4060             mp->pr_mflags = 0;
4061             if (prot & PROT_READ)
4062                 mp->pr_mflags |= MA_READ;
4063             if (prot & PROT_WRITE)
4064                 mp->pr_mflags |= MA_WRITE;
4065             if (prot & PROT_EXEC)
4066                 mp->pr_mflags |= MA_EXEC;
4067             if (segop_gettype(seg, saddr) & MAP_SHARED)
4068                 if (SEGOP_GETTYPE(seg, saddr) & MAP_SHARED)
4069                     mp->pr_mflags |= MA_SHARED;
4070             if (segop_gettype(seg, saddr) & MAP_NORESERVE)
4071                 if (SEGOP_GETTYPE(seg, saddr) & MAP_NORESERVE)
4072                     mp->pr_mflags |= MA_NORESERVE;
4073             if (seg->s_ops == &segspt_shmops ||
4074                 (seg->s_ops == &segvn_ops &&
4075                  (segop_getvp(seg, saddr, &vp) != 0 ||
4076                   (SEGOP_GETVP(seg, saddr, &vp) != 0 ||
4077                    vp == NULL)))
4078                 mp->pr_mflags |= MA_ANON;
4079             if (seg == brkseg)
4080                 mp->pr_mflags |= MA_BREAK;
4081             else if (seg == stkseg)
4082                 mp->pr_mflags |= MA_STACK;
4083             if (seg->s_ops == &segspt_shmops)
4084                 mp->pr_mflags |= MA_ISM | MA_SHM;

4086             mp->pr_pagesize = PAGE_SIZE;
4087             if (psz == -1) {
4088                 mp->pr_hatpagesize = 0;
4089             } else {
4090                 mp->pr_hatpagesize = psz;
4091             }

```

```

4089         /*
4090          * Manufacture a filename for the "object" dir.
4091          */
4092         mp->pr_dev = PRNODEV;
4093         vattr.va_mask = AT_FSID|AT_NODEID;
4094         if (seg->s_ops == &segvn_ops &&
4095             segop_getvp(seg, saddr, &vp) == 0 &&
4096             SEGOP_GETVP(seg, saddr, &vp) == 0 &&
4097             vp != NULL && vp->v_type == VREG &&
4098             VOP_GETATTR(vp, &vattr, 0, CRED(),
4099             NULL) == 0) {
4100             mp->pr_dev = vattr.va_fsid;
4101             mp->pr_ino = vattr.va_nodeid;
4102             if (vp == p->p_exec)
4103                 (void) strcpy(mp->pr_mapname,
4104                 "a.out");
4105             else
4106                 pr_object_name(mp->pr_mapname,
4107                 vp, &vattr);
4108         }

4109     /*
4110      * Get the SysV shared memory id, if any.
4111      */
4112     if ((mp->pr_mflags & MA_SHARED) &&
4113         p->p_segacct && (mp->pr_shmid = shmgetid(p,
4114         seg->s_base) != SHMID_NONE) {
4115         if (mp->pr_shmid == SHMID_FREE)
4116             mp->pr_shmid = -1;

4118         mp->pr_mflags |= MA_SHM;
4119     } else {
4120         mp->pr_shmid = -1;
4121     }

4123     npages = ((uintptr_t)(naddr - saddr)) >>
4124     PAGE_SHIFT;
4125     parr = kmem_zalloc(npages, KM_SLEEP);

4127     (void) segop_incore(seg, saddr, naddr - saddr,
4128     parr);
4129     SEGOP_INCORE(seg, saddr, naddr - saddr, parr);

4130     for (pagenum = 0; pagenum < npages; pagenum++) {
4131         if (parr[pagenum] & SEG_PAGE_INCORE)
4132             mp->pr_rss++;
4133         if (parr[pagenum] & SEG_PAGE_ANON)
4134             mp->pr_anon++;
4135         if (parr[pagenum] & SEG_PAGE_LOCKED)
4136             mp->pr_locked++;
4137     }
4138     kmem_free(parr, npages);
4139     }
4140     }
4141     ASSERT(tmp == NULL);
4142     } while ((seg = AS_SEGNEXT(as, seg)) != NULL);

4144     return (0);
4145 }

unchanged_portion_omitted

4181 #ifdef _SYSCALL32_IMPL
4182 /*
4183  * Return an array of structures with HAT memory map information.
4184  * We allocate here; the caller must deallocate.

```

```

4185 */
4186 int
4187 prgetxmap32(proc_t *p, list_t *iolhead)
4188 {
4189     struct as *as = p->p_as;
4190     prxmap32_t *mp;
4191     struct seg *seg;
4192     struct seg *brkseg, *stkseg;
4193     struct vnode *vp;
4194     struct vattr vattr;
4195     uint_t prot;

4197     ASSERT(as != &kas && AS_WRITE_HELD(as, &as->a_lock));

4199     /*
4200      * Request an initial buffer size that doesn't waste memory
4201      * if the address space has only a small number of segments.
4202      */
4203     pr_iol_initlist(iolhead, sizeof (*mp), avl_numnodes(&as->a_segtree));

4205     if ((seg = AS_SEGFIRST(as)) == NULL)
4206         return (0);

4208     brkseg = break_seg(p);
4209     stkseg = as_segat(as, prgetstackbase(p));

4211     do {
4212         caddr_t eaddr = seg->s_base + pr_getsegsz(seg, 0);
4213         caddr_t saddr, naddr, baddr;
4214         void *tmp = NULL;
4215         ssize_t psz;
4216         char *parr;
4217         uint64_t npages;
4218         uint64_t pagenum;

4220         /*
4221          * Segment loop part one: iterate from the base of the segment
4222          * to its end, pausing at each address boundary (baddr) between
4223          * ranges that have different virtual memory protections.
4224          */
4225         for (saddr = seg->s_base; saddr < eaddr; saddr = baddr) {
4226             prot = pr_getprot(seg, 0, &tmp, &saddr, &baddr, eaddr);
4227             ASSERT(baddr >= saddr && baddr <= eaddr);

4229             /*
4230              * Segment loop part two: iterate from the current
4231              * position to the end of the protection boundary,
4232              * pausing at each address boundary (naddr) between
4233              * ranges that have different underlying page sizes.
4234              */
4235             for (; saddr < baddr; saddr = naddr) {
4236                 psz = pr_getpagesize(seg, saddr, &naddr, baddr);
4237                 ASSERT(naddr >= saddr && naddr <= baddr);

4239                 mp = pr_iol_newbuf(iolhead, sizeof (*mp));

4241                 mp->pr_vaddr = (caddr32_t)(uintptr_t)saddr;
4242                 mp->pr_size = (size32_t)(naddr - saddr);
4243                 mp->pr_offset = segop_getoffset(seg, saddr);
4244                 mp->pr_offset = SEGOP_GETOFFSET(seg, saddr);
4245                 mp->pr_mflags = 0;
4246                 if (prot & PROT_READ)
4247                     mp->pr_mflags |= MA_READ;
4247                 if (prot & PROT_WRITE)
4248                     mp->pr_mflags |= MA_WRITE;
4248                 if (prot & PROT_EXEC)
4249                     mp->pr_mflags |= MA_EXEC;

```

```

4250         mp->pr_mflags |= MA_EXEC;
4251         if (segop_gettype(seg, saddr) & MAP_SHARED)
4252             if (SEGOP_GETTYPE(seg, saddr) & MAP_SHARED)
4253                 mp->pr_mflags |= MA_SHARED;
4253         if (segop_gettype(seg, saddr) & MAP_NORESERVE)
4254             if (SEGOP_GETTYPE(seg, saddr) & MAP_NORESERVE)
4255                 mp->pr_mflags |= MA_NORESERVE;
4256         if (seg->s_ops == &segspt_shmops ||
4257             (seg->s_ops == &segvn_ops &&
4258              (segop_getvp(seg, saddr, &vp) != 0 ||
4259               (SEGOP_GETVP(seg, saddr, &vp) != 0 ||
4260                vp == NULL)))
4261             mp->pr_mflags |= MA_ANON;
4262         if (seg == brkseg)
4263             mp->pr_mflags |= MA_BREAK;
4264         else if (seg == stkseg)
4265             mp->pr_mflags |= MA_STACK;
4266         if (seg->s_ops == &segspt_shmops)
4267             mp->pr_mflags |= MA_ISM | MA_SHM;

4267     mp->pr_pagesize = PAGESIZE;
4268     if (psz == -1) {
4269         mp->pr_hatpagesize = 0;
4270     } else {
4271         mp->pr_hatpagesize = psz;
4272     }

4274     /*
4275      * Manufacture a filename for the "object" dir.
4276      */
4277     mp->pr_dev = PRNODEV32;
4278     vattr.va_mask = AT_FSID|AT_NODEID;
4279     if (seg->s_ops == &segvn_ops &&
4280         segop_getvp(seg, saddr, &vp) == 0 &&
4281         SEGOP_GETVP(seg, saddr, &vp) == 0 &&
4282         vp != NULL && vp->v_type == VREG &&
4283         VOP_GETATTR(vp, &vattr, 0, CRED(),
4284                     NULL) == 0) {
4285         (void) cmlpdev(&mp->pr_dev,
4286                      vattr.va_fsid);
4287         mp->pr_ino = vattr.va_nodeid;
4288         if (vp == p->p_exec)
4289             (void) strcpy(mp->pr_mapname,
4290                          "a.out");
4291     } else
4292         pr_object_name(mp->pr_mapname,
4293                      vp, &vattr);

4295     /*
4296      * Get the SysV shared memory id, if any.
4297      */
4298     if ((mp->pr_mflags & MA_SHARED) &&
4299         p->p_segacct && (mp->pr_shmid == shmgetid(p,
4300         seg->s_base)) != SHMID_NONE) {
4301         if (mp->pr_shmid == SHMID_FREE)
4302             mp->pr_shmid = -1;

4304         mp->pr_mflags |= MA_SHM;
4305     } else {
4306         mp->pr_shmid = -1;
4307     }

4309     npages = ((uintptr_t)(naddr - saddr)) >>
4310             PAGESHIFT;
4311     parr = kmem_zalloc(npages, KM_SLEEP);

```

```
4313     (void) segop_incore(seg, saddr, naddr - saddr,
4314                       parr);
4312     SEGOP_INCORE(seg, saddr, naddr - saddr, parr);

4316     for (pagenum = 0; pagenum < npages; pagenum++) {
4317         if (parr[pagenum] & SEG_PAGE_INCORE)
4318             mp->pr_rss++;
4319         if (parr[pagenum] & SEG_PAGE_ANON)
4320             mp->pr_anon++;
4321         if (parr[pagenum] & SEG_PAGE_LOCKED)
4322             mp->pr_locked++;
4323     }
4324     kmem_free(parr, npages);
4325 }
4326 }
4327     ASSERT(tmp == NULL);
4328 } while ((seg = AS_SEGNEXT(as, seg)) != NULL);

4330     return (0);
4331 }
_____unchanged_portion_omitted_
```



```
*****
```

```
142974 Tue Nov 24 09:34:45 2015
```

```
new/usr/src/uts/common/fs/proc/prvnops.c
```

```
patch lower-case-segops
```

```
*****
```

```
_____ unchanged portion omitted _____
```

```
3663 static vnode_t *
3664 pr_lookup_objctdir(vnode_t *dp, char *comp)
3665 {
3666     prnode_t *dnp = VTOP(dp);
3667     prnode_t *pnp;
3668     proc_t *p;
3669     struct seg *seg;
3670     struct as *as;
3671     vnode_t *vp;
3672     vattr_t vattr;
3673
3674     ASSERT(dnp->pr_type == PR_OBJECTDIR);
3675
3676     pnp = prgetnode(dp, PR_OBJECT);
3677
3678     if (prlock(dnp, ZNO) != 0) {
3679         prfreenode(pnp);
3680         return (NULL);
3681     }
3682     p = dnp->pr_common->prc_proc;
3683     if ((p->p_flag & SSYS) || (as = p->p_as) == &kas) {
3684         prunlock(dnp);
3685         prfreenode(pnp);
3686         return (NULL);
3687     }
3688
3689     /*
3690     * We drop p_lock before grabbing the address space lock
3691     * in order to avoid a deadlock with the clock thread.
3692     * The process will not disappear and its address space
3693     * will not change because it is marked P_PR_LOCK.
3694     */
3695     mutex_exit(&p->p_lock);
3696     AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
3697     if ((seg = AS_SEGFIRST(as)) == NULL) {
3698         vp = NULL;
3699         goto out;
3700     }
3701     if (strcmp(comp, "a.out") == 0) {
3702         vp = p->p_exec;
3703         goto out;
3704     }
3705     do {
3706         /*
3707         * Manufacture a filename for the "object" directory.
3708         */
3709         vattr.va_mask = AT_FSID|AT_NODEID;
3710         if (seg->s_ops == &segvn_ops &&
3711             segop_getvp(seg, seg->s_base, &vp) == 0 &&
3712             SEGOP_GETVP(seg, seg->s_base, &vp) == 0 &&
3713             vp != NULL && vp->v_type == VREG &&
3714             VOP_GETATTR(vp, &vattr, 0, CRED(), NULL) == 0) {
3715             char name[64];
3716
3717             if (vp == p->p_exec) /* "a.out" */
3718                 continue;
3719             pr_object_name(name, vp, &vattr);
3720             if (strcmp(name, comp) == 0)
3721                 goto out;
3722         }
3723     } while (seg = seg->segv);
3724     prfreenode(pnp);
3725     return (NULL);
3726 }
```

```
3721     }
3722     } while ((seg = AS_SEGNEXT(as, seg)) != NULL);
3723
3724     vp = NULL;
3725 out:
3726     if (vp != NULL) {
3727         VN_HOLD(vp);
3728     }
3729     AS_LOCK_EXIT(as, &as->a_lock);
3730     mutex_enter(&p->p_lock);
3731     prunlock(dnp);
3732
3733     if (vp == NULL)
3734         prfreenode(pnp);
3735     else {
3736         /*
3737         * Fill in the prnode so future references will
3738         * be able to find the underlying object's vnode.
3739         * Don't link this prnode into the list of all
3740         * prnodes for the process; this is a one-use node.
3741         * Its use is entirely to catch and fail opens for writing.
3742         */
3743         pnp->pr_realvp = vp;
3744         vp = PTOV(pnp);
3745     }
3746
3747     return (vp);
3748 }
3749 _____ unchanged portion omitted _____
3750
3751 static vnode_t *
3752 pr_lookup_pathdir(vnode_t *dp, char *comp)
3753 {
3754     prnode_t *dnp = VTOP(dp);
3755     prnode_t *pnp;
3756     vnode_t *vp = NULL;
3757     proc_t *p;
3758     uint_t fd, flags = 0;
3759     int c;
3760     uf_entry_t *ufp;
3761     uf_info_t *fip;
3762     enum { NAME_FD, NAME_OBJECT, NAME_ROOT, NAME_CWD, NAME_UNKNOWN } type;
3763     char *tmp;
3764     int idx;
3765     struct seg *seg;
3766     struct as *as = NULL;
3767     vattr_t vattr;
3768
3769     ASSERT(dnp->pr_type == PR_PATHDIR);
3770
3771     /*
3772     * First, check if this is a numeric entry, in which case we have a
3773     * file descriptor.
3774     */
3775     fd = 0;
3776     type = NAME_FD;
3777     tmp = comp;
3778     while ((c = *tmp++) != '\0') {
3779         int ofd;
3780         if (c < '0' || c > '9') {
3781             type = NAME_UNKNOWN;
3782             break;
3783         }
3784         ofd = fd;
3785         fd = 10*fd + c - '0';
3786         if (ofd/10 != ofd) { /* integer overflow */
```

```

4087         type = NAME_UNKNOWN;
4088         break;
4089     }
4090 }

4092 /*
4093  * Next, see if it is one of the special values {root, cwd}.
4094  */
4095 if (type == NAME_UNKNOWN) {
4096     if (strcmp(comp, "root") == 0)
4097         type = NAME_ROOT;
4098     else if (strcmp(comp, "cwd") == 0)
4099         type = NAME_CWD;
4100 }

4102 /*
4103  * Grab the necessary data from the process
4104  */
4105 if (prlock(dpnp, ZNO) != 0)
4106     return (NULL);
4107 p = dpnp->pr_common->prc_proc;

4109 fip = P_FINFO(p);

4111 switch (type) {
4112 case NAME_ROOT:
4113     if ((vp = PTOU(p)->u_rdir) == NULL)
4114         vp = p->p_zone->zone_rootvp;
4115     VN_HOLD(vp);
4116     break;
4117 case NAME_CWD:
4118     vp = PTOU(p)->u_cdir;
4119     VN_HOLD(vp);
4120     break;
4121 default:
4122     if ((p->p_flag & SSYS) || (as = p->p_as) == &kas) {
4123         prunlock(dpnp);
4124         return (NULL);
4125     }
4126 }
4127 mutex_exit(&p->p_lock);

4129 /*
4130  * Determine if this is an object entry
4131  */
4132 if (type == NAME_UNKNOWN) {
4133     /*
4134      * Start with the inode index immediately after the number of
4135      * files.
4136      */
4137     mutex_enter(&fip->fi_lock);
4138     idx = fip->fi_nfiles + 4;
4139     mutex_exit(&fip->fi_lock);

4141     if (strcmp(comp, "a.out") == 0) {
4142         if (p->p_execdir != NULL) {
4143             vp = p->p_execdir;
4144             VN_HOLD(vp);
4145             type = NAME_OBJECT;
4146             flags |= PR_AOUT;
4147         } else {
4148             vp = p->p_exec;
4149             VN_HOLD(vp);
4150             type = NAME_OBJECT;
4151         }
4152     } else {

```

```

4153     AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
4154     if ((seg = AS_SEGFIRST(as)) != NULL) {
4155         do {
4156             /*
4157              * Manufacture a filename for the
4158              * "object" directory.
4159              */
4160             vattr.va_mask = AT_FSID|AT_NODEID;
4161             if (seg->s_ops == &segvn_ops &&
4162                 segop_getvp(seg, seg->s_base, &vp)
4163                 SEGOP_GETVP(seg, seg->s_base, &vp)
4164                 == 0 &&
4165                 vp != NULL && vp->v_type == VREG &&
4166                 VOP_GETATTR(vp, &vattr, 0, CRED(),
4167                     NULL) == 0) {
4168                 char name[64];
4169
4170                 if (vp == p->p_exec)
4171                     continue;
4172                 idx++;
4173                 pr_object_name(name, vp,
4174                     &vattr);
4175                 if (strcmp(name, comp) == 0)
4176                     break;
4177             } while ((seg = AS_SEGNEXT(as, seg)) != NULL);
4178         }
4179     }
4180     if (seg == NULL) {
4181         vp = NULL;
4182     } else {
4183         VN_HOLD(vp);
4184         type = NAME_OBJECT;
4185     }
4187     AS_LOCK_EXIT(as, &as->a_lock);
4188 }
4189

4192 switch (type) {
4193 case NAME_FD:
4194     mutex_enter(&fip->fi_lock);
4195     if (fd < fip->fi_nfiles) {
4196         UF_ENTER(ufp, fip, fd);
4197         if (ufp->uf_file != NULL) {
4198             vp = ufp->uf_file->f_vnode;
4199             VN_HOLD(vp);
4200         }
4201         UF_EXIT(ufp);
4202     }
4203     mutex_exit(&fip->fi_lock);
4204     idx = fd + 4;
4205     break;
4206 case NAME_ROOT:
4207     idx = 2;
4208     break;
4209 case NAME_CWD:
4210     idx = 3;
4211     break;
4212 case NAME_OBJECT:
4213 case NAME_UNKNOWN:
4214     /* Nothing to do */
4215     break;
4216 }

```

```

4218     mutex_enter(&p->p_lock);
4219     prunlock(dpnp);

4221     if (vp != NULL) {
4222         pnp = prgetnode(dp, PR_PATH);

4224         pnp->pr_flags |= flags;
4225         pnp->pr_common = dpnp->pr_common;
4226         pnp->pr_pcommon = dpnp->pr_pcommon;
4227         pnp->pr_realvp = vp;
4228         pnp->pr_parent = dp;          /* needed for prlookup */
4229         pnp->pr_ino = pmkino(idx, dpnp->pr_common->prc_slot, PR_PATH);
4230         VN_HOLD(dp);
4231         vp = PTOV(pnp);
4232         vp->v_type = VLNK;
4233     }

4235     return (vp);
4236 }

```

unchanged portion omitted

```

4829 static void
4830 rebuild_objdir(struct as *as)
4831 {
4832     struct seg *seg;
4833     vnode_t *vp;
4834     vattn_t vattn;
4835     vnode_t **dir;
4836     ulong_t nalloc;
4837     ulong_t nentries;
4838     int i, j;
4839     ulong_t nold, nnew;

4841     ASSERT(AS_WRITE_HELD(as, &as->a_lock));

4843     if (as->a_updatedir == 0 && as->a_objectdir != NULL)
4844         return;
4845     as->a_updatedir = 0;

4847     if ((nalloc = avl_numnodes(&as->a_segtree)) == 0 ||
4848         (seg = AS_SEGFIRST(as)) == NULL) /* can't happen? */
4849         return;

4851     /*
4852      * Allocate space for the new object directory.
4853      * (This is usually about two times too many entries.)
4854      */
4855     nalloc = (nalloc + 0xf) & ~0xf; /* multiple of 16 */
4856     dir = kmem_zalloc(nalloc * sizeof (vnode_t *), KM_SLEEP);

4858     /* fill in the new directory with desired entries */
4859     nentries = 0;
4860     do {
4861         vattn.va_mask = AT_FSID|AT_NODEID;
4862         if (seg->s_ops == &segvn_ops &&
4863             segop_getvp(seg, seg->s_base, &vp) == 0 &&
4864             SEGOP_GETVP(seg, seg->s_base, &vp) == 0 &&
4865             vp != NULL && vp->v_type == VREG &&
4866             VOP_GETATTR(vp, &vattn, 0, CRED(), NULL) == 0) {
4867             for (i = 0; i < nentries; i++)
4868                 if (vp == dir[i])
4869                     break;
4870             if (i == nentries) {
4871                 ASSERT(nentries < nalloc);
4872                 dir[nentries++] = vp;
4873             }
4874         }

```

```

4873     }
4874     } while ((seg = AS_SEGNEXT(as, seg)) != NULL);

4876     if (as->a_objectdir == NULL) { /* first time */
4877         as->a_objectdir = dir;
4878         as->a_sizedir = nalloc;
4879         return;
4880     }

4882     /*
4883      * Null out all of the defunct entries in the old directory.
4884      */
4885     nold = 0;
4886     nnew = nentries;
4887     for (i = 0; i < as->a_sizedir; i++) {
4888         if ((vp = as->a_objectdir[i]) != NULL) {
4889             for (j = 0; j < nentries; j++) {
4890                 if (vp == dir[j]) {
4891                     dir[j] = NULL;
4892                     nnew--;
4893                     break;
4894                 }
4895             }
4896             if (j == nentries)
4897                 as->a_objectdir[i] = NULL;
4898             else
4899                 nold++;
4900         }
4901     }

4903     if (nold + nnew > as->a_sizedir) {
4904         /*
4905          * Reallocate the old directory to have enough
4906          * space for the old and new entries combined.
4907          * Round up to the next multiple of 16.
4908          */
4909         ulong_t newsize = (nold + nnew + 0xf) & ~0xf;
4910         vnode_t **newdir = kmem_zalloc(newsize * sizeof (vnode_t *),
4911             KM_SLEEP);
4912         bcopy(as->a_objectdir, newdir,
4913             as->a_sizedir * sizeof (vnode_t *));
4914         kmem_free(as->a_objectdir, as->a_sizedir * sizeof (vnode_t *));
4915         as->a_objectdir = newdir;
4916         as->a_sizedir = newsize;
4917     }

4919     /*
4920      * Move all new entries to the old directory and
4921      * deallocate the space used by the new directory.
4922      */
4923     if (nnew) {
4924         for (i = 0, j = 0; i < nentries; i++) {
4925             if ((vp = dir[i]) == NULL)
4926                 continue;
4927             for (; j < as->a_sizedir; j++) {
4928                 if (as->a_objectdir[j] != NULL)
4929                     continue;
4930                 as->a_objectdir[j++] = vp;
4931                 break;
4932             }
4933         }
4934     }
4935     kmem_free(dir, nalloc * sizeof (vnode_t *));
4936 }

```

unchanged portion omitted

new/usr/src/uts/common/io/mem.c

1

23668 Tue Nov 24 09:34:45 2015

new/usr/src/uts/common/io/mem.c

patch lower-case-segops

_____unchanged_portion_omitted_____

```
285 static int
286 mmpagelock(struct as *as, caddr_t va)
287 {
288     struct seg *seg;
289     int i;

291     AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
292     seg = as_segat(as, va);
293     i = (seg != NULL)? segop_capable(seg, S_CAPABILITY_NOMINFLT) : 0;
293     i = (seg != NULL)? SEGOP_CAPABLE(seg, S_CAPABILITY_NOMINFLT) : 0;
294     AS_LOCK_EXIT(as, &as->a_lock);

296     return (i);
297 }
_____unchanged_portion_omitted_____
```

new/usr/src/uts/common/os/dumpsubr.c

1

80155 Tue Nov 24 09:34:46 2015

new/usr/src/uts/common/os/dumpsubr.c

patch lower-case-segops

_____unchanged_portion_omitted_____

```
1401 /*
1402  * Dump the <as, va, pfn> information for a given address space.
1403  * segop_dump() will call dump_addpage() for each page in the segment.
1403  * SEGOP_DUMP() will call dump_addpage() for each page in the segment.
1404  */
1405 static void
1406 dump_as(struct as *as)
1407 {
1408     struct seg *seg;
1409
1410     AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
1411     for (seg = AS_SEGFIRST(as); seg; seg = AS_SEGNEXT(as, seg)) {
1412         if (seg->s_as != as)
1413             break;
1414         if (seg->s_ops == NULL)
1415             continue;
1416         segop_dump(seg);
1416         SEGOP_DUMP(seg);
1417     }
1418     AS_LOCK_EXIT(as, &as->a_lock);
1419
1420     if (seg != NULL)
1421         cmn_err(CE_WARN, "invalid segment %p in address space %p",
1422             (void *)seg, (void *)as);
1423 }
```

_____unchanged_portion_omitted_____

```

*****
52333 Tue Nov 24 09:34:46 2015
new/usr/src/uts/common/os/exec.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

1144 /*
1145  * Map a section of an executable file into the user's
1146  * address space.
1147  */
1148 int
1149 execmap(struct vnode *vp, caddr_t addr, size_t len, size_t zfodlen,
1150         off_t offset, int prot, int page, uint_t szc)
1151 {
1152     int error = 0;
1153     off_t oldoffset;
1154     caddr_t zfodbase, oldaddr;
1155     size_t end, oldlen;
1156     size_t zfoddiff;
1157     label_t ljb;
1158     proc_t *p = ttoproc(curthread);

1160     oldaddr = addr;
1161     addr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);
1162     if (len) {
1163         oldlen = len;
1164         len += ((size_t)oldaddr - (size_t)addr);
1165         oldoffset = offset;
1166         offset = (off_t)((uintptr_t)offset & PAGEMASK);
1167         if (page) {
1168             spgcnt_t prefltmem, availm, npages;
1169             int preread;
1170             uint_t mflag = MAP_PRIVATE | MAP_FIXED;

1172             if ((prot & (PROT_WRITE | PROT_EXEC)) == PROT_EXEC) {
1173                 mflag |= MAP_TEXT;
1174             } else {
1175                 mflag |= MAP_INITDATA;
1176             }

1178             if (valid_usr_range(addr, len, prot, p->p_as,
1179                                 p->p_as->a_userlimit) != RANGE_OKAY) {
1180                 error = ENOMEM;
1181                 goto bad;
1182             }
1183             if (error = VOP_MAP(vp, (offset_t)offset,
1184                                 p->p_as, &addr, len, prot, PROT_ALL,
1185                                 mflag, CRED(), NULL))
1186                 goto bad;

1188             /*
1189              * If the segment can fit, then we prefault
1190              * the entire segment in. This is based on the
1191              * model that says the best working set of a
1192              * small program is all of its pages.
1193              */
1194             npages = (spgcnt_t)btopr(len);
1195             prefltmem = freemem - desfree;
1196             preread =
1197                 (npages < prefltmem && len < PGTHRESH) ? 1 : 0;

1199             /*
1200              * If we aren't prefaulting the segment,
1201              * increment "deficit", if necessary to ensure
1202              * that pages will become available when this

```

```

1203         * process starts executing.
1204         */
1205         availm = freemem - lotsfree;
1206         if (preread == 0 && npages > availm &&
1207             deficit < lotsfree) {
1208             deficit += MIN((pgcnt_t)(npages - availm),
1209                           lotsfree - deficit);
1210         }

1212         if (preread) {
1213             TRACE_2(TR_FAC_PROC, TR_EXECMAP_PREREAD,
1214                   "execmap preread:freemem %d size %lu",
1215                   freemem, len);
1216             (void) as_fault(p->p_as->a_hat, p->p_as,
1217                             (caddr_t)addr, len, F_INVALID, S_READ);
1218         }
1219     } else {
1220         if (valid_usr_range(addr, len, prot, p->p_as,
1221                             p->p_as->a_userlimit) != RANGE_OKAY) {
1222             error = ENOMEM;
1223             goto bad;
1224         }

1226         if (error = as_map(p->p_as, addr, len,
1227                             segvn_create, zfod_argsp))
1228             goto bad;

1229         /*
1230          * Read in the segment in one big chunk.
1231          */
1232         if (error = vn_rdwr(UIO_READ, vp, (caddr_t)oldaddr,
1233                             oldlen, (offset_t)oldoffset, UIO_USERSPACE, 0,
1234                             (rlim64_t)0, CRED(), (ssize_t *)0))
1235             goto bad;

1236         /*
1237          * Now set protections.
1238          */
1239         if (prot != PROT_ZFOD) {
1240             (void) as_setprot(p->p_as, (caddr_t)addr,
1241                               len, prot);
1242         }
1243     }
1244 }

1246 if (zfodlen) {
1247     struct as *as = curproc->p_as;
1248     struct seg *seg;
1249     uint_t zprot = 0;

1251     end = (size_t)addr + len;
1252     zfodbase = (caddr_t)roundup(end, PAGESIZE);
1253     zfoddiff = (uintptr_t)zfodbase - end;
1254     if (zfoddiff) {
1255         /*
1256          * Before we go to zero the remaining space on the last
1257          * page, make sure we have write permission.
1258          *
1259          * Normal illumos binaries don't even hit the case
1260          * where we have to change permission on the last page
1261          * since their protection is typically either
1262          * PROT_USER | PROT_WRITE | PROT_READ
1263          * or
1264          * PROT_ZFOD (same as PROT_ALL).
1265          *
1266          * We need to be careful how we zero-fill the last page
1267          * if the segment protection does not include
1268          * PROT_WRITE. Using as_setprot() can cause the VM

```

```

1269      * segment code to call segvn_vpage(), which must
1270      * allocate a page struct for each page in the segment.
1271      * If we have a very large segment, this may fail, so
1272      * we have to check for that, even though we ignore
1273      * other return values from as_setprot.
1274      */

1276      AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
1277      seg = as_segat(curproc->p_as, (caddr_t)end);
1278      if (seg != NULL)
1279          (void) segop_getprot(seg, (caddr_t)end,
1280                              zfoddiff - 1, &zprot);
1281          SEGOP_GETPROT(seg, (caddr_t)end, zfoddiff - 1,
1282                      &zprot);
1283          AS_LOCK_EXIT(as, &as->a_lock);

1284      if (seg != NULL && (zprot & PROT_WRITE) == 0) {
1285          if (as_setprot(as, (caddr_t)end, zfoddiff - 1,
1286                      zprot | PROT_WRITE) == ENOMEM) {
1287              error = ENOMEM;
1288              goto bad;
1289          }
1290      }

1291      if (on_fault(&ljb)) {
1292          no_fault();
1293          if (seg != NULL && (zprot & PROT_WRITE) == 0)
1294              (void) as_setprot(as, (caddr_t)end,
1295                              zfoddiff - 1, zprot);
1296          error = EFAULT;
1297          goto bad;
1298      }
1299      uzero((void *)end, zfoddiff);
1300      no_fault();
1301      if (seg != NULL && (zprot & PROT_WRITE) == 0)
1302          (void) as_setprot(as, (caddr_t)end,
1303                          zfoddiff - 1, zprot);
1304      }
1305      if (zfodlen > zfoddiff) {
1306          struct segvn_crargs crargs =
1307              SEGVN_ZFOD_ARGS(PROT_ZFOD, PROT_ALL);

1309          zfodlen -= zfoddiff;
1310          if (valid_usr_range(zfodbase, zfodlen, prot, p->p_as,
1311                          p->p_as->a_userlimit) != RANGE_OKAY) {
1312              error = ENOMEM;
1313              goto bad;
1314          }
1315          if (szc > 0) {
1316              /*
1317               * ASSERT alignment because the mapelfexec()
1318               * caller for the szc > 0 case extended zfod
1319               * so it's end is pgsz aligned.
1320               */
1321              size_t pgsz = page_get_pagesize(szc);
1322              ASSERT(IS_P2ALIGNED(zfodbase + zfodlen, pgsz));

1324              if (IS_P2ALIGNED(zfodbase, pgsz)) {
1325                  crargs.szc = szc;
1326              } else {
1327                  crargs.szc = AS_MAP_HEAP;
1328              }
1329          } else {
1330              crargs.szc = AS_MAP_NO_LPOOB;
1331          }
1332      }
1333      if (error = as_map(p->p_as, (caddr_t)zfodbase,

```

```

1333          zfodlen, segvn_create, &crargs))
1334          goto bad;
1335          if (prot != PROT_ZFOD) {
1336              (void) as_setprot(p->p_as, (caddr_t)zfodbase,
1337                              zfodlen, prot);
1338          }
1339      }
1340      }
1341      return (0);
1342 bad:
1343      return (error);
1344 }

```

unchanged portion omitted

new/usr/src/uts/common/os/lgrp.c

1

119448 Tue Nov 24 09:34:46 2015

new/usr/src/uts/common/os/lgrp.c

patch lower-case-segops

unchanged_portion_omitted_

```
3498 /*
3499  * Get memory allocation policy for this segment
3500  */
3501 lgrp_mem_policy_info_t *
3502 lgrp_mem_policy_get(struct seg *seg, caddr_t vaddr)
3503 {
3504     lgrp_mem_policy_info_t *policy_info;
3505     extern struct seg_ops  segspt_ops;
3506     extern struct seg_ops  segspt_shmops;
3507
3508     /*
3509      * This is for binary compatibility to protect against third party
3510      * segment drivers which haven't recompiled to allow for
3511      * segop_getpolicy()
3512      * SEGOP_GETPOLICY()
3513      */
3514     if (seg->s_ops != &segvn_ops && seg->s_ops != &segspt_ops &&
3515         seg->s_ops != &segspt_shmops)
3516         return (NULL);
3517
3518     policy_info = NULL;
3519     if (seg->s_ops->getpolicy != NULL)
3520         policy_info = segop_getpolicy(seg, vaddr);
3521         policy_info = SEGOP_GETPOLICY(seg, vaddr);
3522
3523     return (policy_info);
3524 }
unchanged_portion_omitted_
```



```

*****
69060 Tue Nov 24 09:34:46 2015
new/usr/src/uts/common/os/mmapobj.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

1434 /*
1435  * Check the address space to see if the virtual addresses to be used are
1436  * available.  If they are not, return errno for failure.  On success, 0
1437  * will be returned, and the virtual addresses for each mmapobj_result_t
1438  * will be reserved.  Note that a reservation could have earlier been made
1439  * for a given segment via a /dev/null mapping.  If that is the case, then
1440  * we can use that VA space for our mappings.
1441  * Note: this function will only be used for ET_EXEC binaries.
1442  */
1443 int
1444 check_exec_addrs(int loadable, mmapobj_result_t *mrp, caddr_t start_addr)
1445 {
1446     int i;
1447     struct as *as = curproc->p_as;
1448     struct segvn_crargs crargs = SEGVN_ZFOD_ARGS(PROT_ZFOD, PROT_ALL);
1449     int ret;
1450     caddr_t myaddr;
1451     size_t mylen;
1452     struct seg *seg;

1454     /* No need to reserve swap space now since it will be reserved later */
1455     crargs.flags |= MAP_NORESERVE;
1456     as_rangelock(as);
1457     for (i = 0; i < loadable; i++) {

1459         myaddr = start_addr + (size_t)mrp[i].mr_addr;
1460         mylen = mrp[i].mr_msize;

1462         /* See if there is a hole in the as for this range */
1463         if (as_gap(as, mylen, &myaddr, &mylen, 0, NULL) == 0) {
1464             ASSERT(myaddr == start_addr + (size_t)mrp[i].mr_addr);
1465             ASSERT(mylen == mrp[i].mr_msize);

1467 #ifdef DEBUG
1468             if (MR_GET_TYPE(mrp[i].mr_flags) == MR_PADDING) {
1469                 MOBJ_STAT_ADD(exec_padding);
1470             }
1471 #endif
1472             ret = as_map(as, myaddr, mylen, segvn_create, &crargs);
1473             if (ret) {
1474                 as_rangeunlock(as);
1475                 mmapobj_unmap_exec(mrp, i, start_addr);
1476                 return (ret);
1477             }
1478         } else {
1479             /*
1480              * There is a mapping that exists in the range
1481              * so check to see if it was a "reservation"
1482              * from /dev/null.  The mapping is from
1483              * /dev/null if the mapping comes from
1484              * segdev and the type is neither MAP_SHARED
1485              * nor MAP_PRIVATE.
1486              */
1487             AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
1488             seg = as_findseg(as, myaddr, 0);
1489             MOBJ_STAT_ADD(exec_addr_mapped);
1490             if (seg && seg->s_ops == &segdev_ops &&
1491                 ((segop_gettype(seg, myaddr) &
1492                  (SEGOP_GETTYPE(seg, myaddr) &

```

```

1492         (MAP_SHARED | MAP_PRIVATE)) == 0) &&
1493         myaddr >= seg->s_base &&
1494         myaddr + mylen <=
1495         seg->s_base + seg->s_size) {
1496             MOBJ_STAT_ADD(exec_addr_devnull);
1497             AS_LOCK_EXIT(as, &as->a_lock);
1498             (void) as_unmap(as, myaddr, mylen);
1499             ret = as_map(as, myaddr, mylen, segvn_create,
1500                 &crargs);
1501             mrp[i].mr_flags |= MR_RESV;
1502             if (ret) {
1503                 as_rangeunlock(as);
1504                 /* Need to remap what we unmapped */
1505                 mmapobj_unmap_exec(mrp, i + 1,
1506                     start_addr);
1507                 return (ret);
1508             }
1509         } else {
1510             AS_LOCK_EXIT(as, &as->a_lock);
1511             as_rangeunlock(as);
1512             mmapobj_unmap_exec(mrp, i, start_addr);
1513             MOBJ_STAT_ADD(exec_addr_in_use);
1514             return (EADDRINUSE);
1515         }
1516     }
1517     as_rangeunlock(as);
1518     return (0);
1519 }
1520 }
_____unchanged_portion_omitted_____

```

```

*****
248841 Tue Nov 24 09:34:46 2015
new/usr/src/uts/common/os/sunddi.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

8219 /*
8220 * A consolidation private function which is essentially equivalent to
8221 * ddi_umem_lock but with the addition of arguments ops_vector and procp.
8222 * A call to as_add_callback is done if DDI_UMEMLOCK_LONGTERM is set, and
8223 * the ops_vector is valid.
8224 *
8225 * Lock the virtual address range in the current process and create a
8226 * ddi_umem_cookie (of type UMEM_LOCKED). This can be used to pass to
8227 * ddi_umem_iosetup to create a buf or do devmap_umem_setup/remap to export
8228 * to user space.
8229 *
8230 * Note: The resource control accounting currently uses a full charge model
8231 * in other words attempts to lock the same/overlapping areas of memory
8232 * will deduct the full size of the buffer from the projects running
8233 * counter for the device locked memory.
8234 *
8235 * addr, size should be PAGESIZE aligned
8236 *
8237 * flags - DDI_UMEMLOCK_READ, DDI_UMEMLOCK_WRITE or both
8238 * identifies whether the locked memory will be read or written or both
8239 * DDI_UMEMLOCK_LONGTERM must be set when the locking will
8240 * be maintained for an indefinitely long period (essentially permanent),
8241 * rather than for what would be required for a typical I/O completion.
8242 * When DDI_UMEMLOCK_LONGTERM is set, umem_lockmemory will return EFAULT
8243 * if the memory pertains to a regular file which is mapped MAP_SHARED.
8244 * This is to prevent a deadlock if a file truncation is attempted after
8245 * after the locking is done.
8246 *
8247 * Returns 0 on success
8248 * EINVAL - for invalid parameters
8249 * EPERM, ENOMEM and other error codes returned by as_pagelock
8250 * ENOMEM - is returned if the current request to lock memory exceeds
8251 * *.max-locked-memory resource control value.
8252 * EFAULT - memory pertains to a regular file mapped shared and
8253 * and DDI_UMEMLOCK_LONGTERM flag is set
8254 * EAGAIN - could not start the ddi_umem_unlock list processing thread
8255 */
8256 int
8257 umem_lockmemory(caddr_t addr, size_t len, int flags, ddi_umem_cookie_t *cookie,
8258                struct umem_callback_ops *ops_vector,
8259                proc_t *procp)
8260 {
8261     int error;
8262     struct ddi_umem_cookie *p;
8263     void (*driver_callback)() = NULL;
8264     struct as *as;
8265     struct seg *seg;
8266     vnode_t *vp;

8268     /* Allow device drivers to not have to reference "curproc" */
8269     if (procp == NULL)
8270         procp = curproc;
8271     as = procp->p_as;
8272     *cookie = NULL; /* in case of any error return */

8274     /* These are the only three valid flags */
8275     if ((flags & ~(DDI_UMEMLOCK_READ | DDI_UMEMLOCK_WRITE |
8276                  DDI_UMEMLOCK_LONGTERM)) != 0)
8277         return (EINVAL);

```

```

8279     /* At least one (can be both) of the two access flags must be set */
8280     if ((flags & (DDI_UMEMLOCK_READ | DDI_UMEMLOCK_WRITE)) == 0)
8281         return (EINVAL);

8283     /* addr and len must be page-aligned */
8284     if (((uintptr_t)addr & PAGEOFFSET) != 0)
8285         return (EINVAL);

8287     if ((len & PAGEOFFSET) != 0)
8288         return (EINVAL);

8290     /*
8291     * For longterm locking a driver callback must be specified; if
8292     * not longterm then a callback is optional.
8293     */
8294     if (ops_vector != NULL) {
8295         if (ops_vector->cbo_umem_callback_version !=
8296             UMEM_CALLBACK_VERSION)
8297             return (EINVAL);
8298         else
8299             driver_callback = ops_vector->cbo_umem_lock_cleanup;
8300     }
8301     if ((driver_callback == NULL) && (flags & DDI_UMEMLOCK_LONGTERM))
8302         return (EINVAL);

8304     /*
8305     * Call i_ddi_umem_unlock_thread_start if necessary. It will
8306     * be called on first ddi_umem_lock or umem_lockmemory call.
8307     */
8308     if (ddi_umem_unlock_thread == NULL)
8309         i_ddi_umem_unlock_thread_start();

8311     /* Allocate memory for the cookie */
8312     p = kmem_zalloc(sizeof (struct ddi_umem_cookie), KM_SLEEP);

8314     /* Convert the flags to seg_rw type */
8315     if (flags & DDI_UMEMLOCK_WRITE) {
8316         p->s_flags = S_WRITE;
8317     } else {
8318         p->s_flags = S_READ;
8319     }

8321     /* Store procp in cookie for later iosetup/unlock */
8322     p->procp = (void *)procp;

8324     /*
8325     * Store the struct as pointer in cookie for later use by
8326     * ddi_umem_unlock. The proc->p_as will be stale if ddi_umem_unlock
8327     * is called after relvm is called.
8328     */
8329     p->asp = as;

8331     /*
8332     * The size field is needed for lockmem accounting.
8333     */
8334     p->size = len;
8335     init_lockedmem_rctl_flag(p);

8337     if (umem_incr_devlockmem(p) != 0) {
8338         /*
8339         * The requested memory cannot be locked
8340         */
8341         kmem_free(p, sizeof (struct ddi_umem_cookie));
8342         *cookie = (ddi_umem_cookie_t)NULL;
8343         return (ENOMEM);

```

```

8344     }
8345
8346     /* Lock the pages corresponding to addr, len in memory */
8347     error = as_pagelock(as, &(p->pparray), addr, len, p->s_flags);
8348     if (error != 0) {
8349         umem_decr_devlockmem(p);
8350         kmem_free(p, sizeof (struct ddi_umem_cookie));
8351         *cookie = (ddi_umem_cookie_t)NULL;
8352         return (error);
8353     }
8354
8355     /*
8356     * For longterm locking the addr must pertain to a seg_vn segment or
8357     * or a seg_spt segment.
8358     * If the segment pertains to a regular file, it cannot be
8359     * mapped MAP_SHARED.
8360     * This is to prevent a deadlock if a file truncation is attempted
8361     * after the locking is done.
8362     * Doing this after as_pagelock guarantees persistence of the as; if
8363     * an unacceptable segment is found, the cleanup includes calling
8364     * as_pageunlock before returning EFAULT.
8365     *
8366     * segdev is allowed here as it is already locked. This allows
8367     * for memory exported by drivers through mmap() (which is already
8368     * locked) to be allowed for LONGTERM.
8369     */
8370     if (flags & DDI_UMEMLOCK_LONGTERM) {
8371         extern struct seg_ops segspt_shmops;
8372         extern struct seg_ops segdev_ops;
8373         AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
8374         for (seg = as_segat(as, addr); ; seg = AS_SEGNEXT(as, seg)) {
8375             if (seg == NULL || seg->s_base > addr + len)
8376                 break;
8377             if (seg->s_ops == &segdev_ops)
8378                 continue;
8379             if (((seg->s_ops != &segvn_ops) &&
8380                 (seg->s_ops != &segspt_shmops)) ||
8381                 ((segop_getvp(seg, addr, &vp) == 0 &&
8382                 ((SEGOP_GETVP(seg, addr, &vp) == 0 &&
8383                 vp != NULL && vp->v_type == VREG) &&
8384                 (segop_gettype(seg, addr) & MAP_SHARED))) {
8385                 (SEGOP_GETTYPE(seg, addr) & MAP_SHARED))) {
8386                 as_pageunlock(as, p->pparray,
8387                     addr, len, p->s_flags);
8388                 AS_LOCK_EXIT(as, &as->a_lock);
8389                 umem_decr_devlockmem(p);
8390                 kmem_free(p, sizeof (struct ddi_umem_cookie));
8391                 *cookie = (ddi_umem_cookie_t)NULL;
8392                 return (EFAULT);
8393             }
8394         }
8395         AS_LOCK_EXIT(as, &as->a_lock);
8396     }
8397
8398     /* Initialize the fields in the ddi_umem_cookie */
8399     p->cvaddr = addr;
8400     p->type = UMEM_LOCKED;
8401     if (driver_callback != NULL) {
8402         /* i_ddi_umem_unlock and umem_lock_undo may need the cookie */
8403         p->cook_refcnt = 2;
8404         p->callbacks = *ops_vector;
8405     } else {
8406         /* only i_ddi_umme_unlock needs the cookie */
8407         p->cook_refcnt = 1;
8408     }

```

```

8409     *cookie = (ddi_umem_cookie_t)p;
8410
8411     /*
8412     * If a driver callback was specified, add an entry to the
8413     * as struct callback list. The as_pagelock above guarantees
8414     * the persistence of as.
8415     */
8416     if (driver_callback) {
8417         error = as_add_callback(as, umem_lock_undo, p, AS_ALL_EVENT,
8418             addr, len, KM_SLEEP);
8419         if (error != 0) {
8420             as_pageunlock(as, p->pparray,
8421                 addr, len, p->s_flags);
8422             umem_decr_devlockmem(p);
8423             kmem_free(p, sizeof (struct ddi_umem_cookie));
8424             *cookie = (ddi_umem_cookie_t)NULL;
8425         }
8426     }
8427     return (error);
8428 }
_____unchanged_portion_omitted_____

```

```

*****
8540 Tue Nov 24 09:34:47 2015
new/usr/src/uts/common/os/urw.c
patch lower-case-segops
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */

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

29 #pragma ident      "%Z%M% %I%      %E% SMI"

29 #include <sys/atomic.h>
30 #include <sys/errno.h>
31 #include <sys/stat.h>
32 #include <sys/modctl.h>
33 #include <sys/conf.h>
34 #include <sys/system.h>
35 #include <sys/ddi.h>
36 #include <sys/sunddi.h>
37 #include <sys/cpuvar.h>
38 #include <sys/kmem.h>
39 #include <sys/strsubr.h>
40 #include <sys/sysmacros.h>
41 #include <sys/frame.h>
42 #include <sys/stack.h>
43 #include <sys/proc.h>
44 #include <sys/priv.h>
45 #include <sys/policy.h>
46 #include <sys/ontrap.h>
47 #include <sys/vmsystem.h>
48 #include <sys/prsystem.h>

50 #include <vm/as.h>
51 #include <vm/seg.h>
52 #include <vm/seg_dev.h>
53 #include <vm/seg_vn.h>
54 #include <vm/seg_spt.h>
55 #include <vm/seg_kmem.h>

57 extern struct seg_ops segdev_ops;      /* needs a header file */
58 extern struct seg_ops segspt_shmops;  /* needs a header file */

```

```

60 static int
61 page_valid(struct seg *seg, caddr_t addr)
62 {
63     struct segvn_data *svd;
64     vnode_t *vp;
65     vattr_t vattr;
66
67     /*
68      * Fail if the page doesn't map to a page in the underlying
69      * mapped file, if an underlying mapped file exists.
70      */
71     vattr.va_mask = AT_SIZE;
72     if (seg->s_ops == &segvn_ops &&
73         segop_getvp(seg, addr, &vp) == 0 &&
74         SEGOP_GETVP(seg, addr, &vp) == 0 &&
75         vp != NULL && vp->v_type == VREG &&
76         VOP_GETATTR(vp, &vattr, 0, CRED(), NULL) == 0) {
77         u_offset_t size = roundup(vattr.va_size, (u_offset_t)PAGESIZE);
78         u_offset_t offset = segop_getoffset(seg, addr);
79         u_offset_t offset = SEGOP_GETOFFSET(seg, addr);
80
81         if (offset >= size)
82             return (0);
83     }
84
85     /*
86      * Fail if this is an ISM shared segment and the address is
87      * not within the real size of the spt segment that backs it.
88      */
89     if (seg->s_ops == &segspt_shmops &&
90         addr >= seg->s_base + spt_realsize(seg))
91         return (0);
92
93     /*
94      * Fail if the segment is mapped from /dev/null.
95      * The key is that the mapping comes from segdev and the
96      * type is neither MAP_SHARED nor MAP_PRIVATE.
97      */
98     if (seg->s_ops == &segdev_ops &&
99         ((segop_gettype(seg, addr) & (MAP_SHARED | MAP_PRIVATE)) == 0) &&
100         ((SEGOP_GETTYPE(seg, addr) & (MAP_SHARED | MAP_PRIVATE)) == 0))
101         return (0);
102
103     /*
104      * Fail if the page is a MAP_NORESERVE page that has
105      * not actually materialized.
106      * We cheat by knowing that segvn is the only segment
107      * driver that supports MAP_NORESERVE.
108      */
109     if (seg->s_ops == &segvn_ops &&
110         (svd = (struct segvn_data *)seg->s_data) != NULL &&
111         (svd->vp == NULL || svd->vp->v_type != VREG) &&
112         (svd->flags & MAP_NORESERVE)) {
113         /*
114          * Guilty knowledge here. We know that
115          * segvn_incore returns more than just the
116          * low-order bit that indicates the page is
117          * actually in memory. If any bits are set,
118          * then there is backing store for the page.
119          */
120         char incore = 0;
121         (void) segop_incore(seg, addr, PAGESIZE, &incore);
122         (void) SEGOP_INCORE(seg, addr, PAGESIZE, &incore);
123         if (incore == 0)
124             return (0);
125     }
126 }

```

```

122     return (1);
123 }
    unchanged_portion_omitted

176 /*
177  * Perform I/O to a given process. This will return EIO if we detect
178  * corrupt memory and ENXIO if there is no such mapped address in the
179  * user process's address space.
180  */
181 static int
182 urw(proc_t *p, int writing, void *buf, size_t len, uintptr_t a)
183 {
184     caddr_t addr = (caddr_t)a;
185     caddr_t page;
186     caddr_t vaddr;
187     struct seg *seg;
188     int error = 0;
189     int err = 0;
190     uint_t prot;
191     uint_t prot_rw = writing ? PROT_WRITE : PROT_READ;
192     int protchanged;
193     on_trap_data_t otd;
194     int retrycnt;
195     struct as *as = p->p_as;
196     enum seg_rw rw;

198     /*
199     * Locate segment containing address of interest.
200     */
201     page = (caddr_t)(uintptr_t)((uintptr_t)addr & PAGEMASK);
202     retrycnt = 0;
203     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
204 retry:
205     if ((seg = as_segat(as, page)) == NULL ||
206         !page_valid(seg, page)) {
207         AS_LOCK_EXIT(as, &as->a_lock);
208         return (ENXIO);
209     }
210     (void) segop_getprot(seg, page, 0, &prot);
211     SEGOP_GETPROT(seg, page, 0, &prot);

212     protchanged = 0;
213     if ((prot & prot_rw) == 0) {
214         protchanged = 1;
215         err = segop_setprot(seg, page, PAGESIZE, prot | prot_rw);
216         err = SEGOP_SETPROT(seg, page, PAGESIZE, prot | prot_rw);

217         if (err == IE_RETRY) {
218             protchanged = 0;
219             ASSERT(retrycnt == 0);
220             retrycnt++;
221             goto retry;
222         }

224         if (err != 0) {
225             AS_LOCK_EXIT(as, &as->a_lock);
226             return (ENXIO);
227         }
228     }

230     /*
231     * segvn may do a copy-on-write for F_SOFTLOCK/S_READ case to break
232     * sharing to avoid a copy on write of a softlocked page by another
233     * thread. But since we locked the address space as a writer no other
234     * thread can cause a copy on write. S_READ_NOCOW is passed as the
235     * access type to tell segvn that it's ok not to do a copy-on-write

```

```

236     * for this SOFTLOCK fault.
237     */
238     if (writing)
239         rw = S_WRITE;
240     else if (seg->s_ops == &segvn_ops)
241         rw = S_READ_NOCOW;
242     else
243         rw = S_READ;

245     if (segop_fault(as->a_hat, seg, page, PAGESIZE, F_SOFTLOCK, rw)) {
246         if (SEGOP_FAULT(as->a_hat, seg, page, PAGESIZE, F_SOFTLOCK, rw)) {
247             if (protchanged)
248                 (void) segop_setprot(seg, page, PAGESIZE, prot);
249                 (void) SEGOP_SETPROT(seg, page, PAGESIZE, prot);
250                 AS_LOCK_EXIT(as, &as->a_lock);
251                 return (ENXIO);
252             }
253             CPU_STATS_ADD_K(vm, softlock, 1);

254     /*
255     * Make sure we're not trying to read or write off the end of the page.
256     */
257     ASSERT(len <= page + PAGESIZE - addr);

258     /*
259     * Map in the locked page, copy to our local buffer,
260     * then map the page out and unlock it.
261     */
262     vaddr = mapin(as, addr, writing);

264     /*
265     * Since we are copying memory on behalf of the user process,
266     * protect against memory error correction faults.
267     */
268     if (!on_trap(&otd, OT_DATA_EC)) {
269         if (seg->s_ops == &segdev_ops) {
270             /*
271             * Device memory can behave strangely; invoke
272             * a segdev-specific copy operation instead.
273             */
274             if (writing) {
275                 if (segdev_copyto(seg, addr, buf, vaddr, len))
276                     error = ENXIO;
277             } else {
278                 if (segdev_copyfrom(seg, addr, vaddr, buf, len))
279                     error = ENXIO;
280             }
281             } else {
282                 if (writing)
283                     bcopy(buf, vaddr, len);
284                 else
285                     bcopy(vaddr, buf, len);
286             }
287         } else {
288             error = EIO;
289         }
290     }
291     no_trap();

292     /*
293     * If we're writing to an executable page, we may need to synchronize
294     * the I$ with the modifications we made through the D$.
295     */
296     if (writing && (prot & PROT_EXEC))
297         sync_icache(vaddr, (uint_t)len);

299     mapout(as, addr, vaddr, writing);

```

```
301     if (rw == S_READ_NOCOW)
302         rw = S_READ;
304     (void) segop_fault(as->a_hat, seg, page, PAGE_SIZE, F_SOFTUNLOCK, rw);
306     (void) SEGOP_FAULT(as->a_hat, seg, page, PAGE_SIZE, F_SOFTUNLOCK, rw);
306     if (protchanged)
307         (void) segop_setprot(seg, page, PAGE_SIZE, prot);
309     (void) SEGOP_SETPROT(seg, page, PAGE_SIZE, prot);
309     AS_LOCK_EXIT(as, &as->a_lock);
311     return (error);
312 }
_____unchanged_portion_omitted_____
```

```

*****
14034 Tue Nov 24 09:34:47 2015
new/usr/src/uts/common/os/vm_subr.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

318 #define MAX_MAPIN_PAGES 8

320 /*
321  * This function temporarily "borrows" user pages for kernel use. If
322  * "cow" is on, it also sets up copy-on-write protection (only feasible
323  * on MAP_PRIVATE segment) on the user mappings, to protect the borrowed
324  * pages from any changes by the user. The caller is responsible for
325  * unlocking and tearing down cow settings when it's done with the pages.
326  * For an example, see kcfree().
327  *
328  * Pages behind [uaddr..uaddr+*lenp] under address space "as" are locked
329  * (shared), and mapped into kernel address range [kaddr..kaddr+*lenp] if
330  * kaddr != -1. On entering this function, cached_ppp contains a list
331  * of pages that are mapped into [kaddr..kaddr+*lenp] already (from a
332  * previous call). Thus if some pages remain behind [uaddr..uaddr+*lenp],
333  * the kernel map won't need to be reloaded again.
334  *
335  * For cow == 1, if the pages are anonymous pages, it also bumps the anon
336  * reference count, and change the user-mapping to read-only. This
337  * scheme should work on all types of segment drivers. But to be safe,
338  * we check against segvn here.
339  *
340  * Since this function is used to emulate copyin() semantic, it checks
341  * to make sure the user-mappings allow "user-read".
342  *
343  * On exit "lenp" contains the number of bytes successfully locked and
344  * mapped in. For the unsuccessful ones, the caller can fall back to
345  * copyin().
346  *
347  * Error return:
348  * ENOTSUP - operation like this is not supported either on this segment
349  * type, or on this platform type.
350  */
351 int
352 cow_mapin(struct as *as, caddr_t uaddr, caddr_t kaddr, struct page **cached_ppp,
353           struct anon **app, size_t *lenp, int cow)
354 {
355     struct      hat *hat;
356     struct seg  *seg;
357     caddr_t     base;
358     page_t     *pp, *ppp[MAX_MAPIN_PAGES];
359     long       i;
360     int        flags;
361     size_t     size, total = *lenp;
362     char       first = 1;
363     faultcode_t res;

365     *lenp = 0;
366     if (cow) {
367         AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
368         seg = as_findseg(as, uaddr, 0);
369         if ((seg == NULL) || ((base = seg->s_base) > uaddr) ||
370             (uaddr + total) > base + seg->s_size) {
371             AS_LOCK_EXIT(as, &as->a_lock);
372             return (EINVAL);
373         }
374     }
375     /*
376     * The COW scheme should work for all segment types.
377     * But to be safe, we check against segvn.

```

```

377     */
378     if (seg->s_ops != &segvn_ops) {
379         AS_LOCK_EXIT(as, &as->a_lock);
380         return (ENOTSUP);
381     } else if ((segop_gettype(seg, uaddr) & MAP_PRIVATE) == 0) {
382     } else if ((SEGOP_GETTYPE(seg, uaddr) & MAP_PRIVATE) == 0) {
383         AS_LOCK_EXIT(as, &as->a_lock);
384         return (ENOTSUP);
385     }
386     hat = as->a_hat;
387     size = total;
388     tryagain:
389     /*
390     * If (cow), hat_softlock will also change the usr protection to RO.
391     * This is the first step toward setting up cow. Before we
392     * bump up an_refcnt, we can't allow any cow-fault on this
393     * address. Otherwise segvn_fault will change the protection back
394     * to RW upon seeing an_refcnt == 1.
395     * The solution is to hold the writer lock on "as".
396     */
397     res = hat_softlock(hat, uaddr, &size, &ppp[0], cow ? HAT_COW : 0);
398     size = total - size;
399     *lenp += size;
400     size = size >> PAGESHIFT;
401     i = 0;
402     while (i < size) {
403         pp = ppp[i];
404         if (cow) {
405             kmutex_t *ahm;
406             /*
407             * Another solution is to hold SE_EXCL on pp, and
408             * disable PROT_WRITE. This also works for MAP_SHARED
409             * segment. The disadvantage is that it locks the
410             * page from being used by anybody else.
411             */
412             ahm = AH_MUTEX(pp->p_vnode, pp->p_offset);
413             mutex_enter(ahm);
414             *app = swap_anon(pp->p_vnode, pp->p_offset);
415             /*
416             * Since we are holding the as lock, this avoids a
417             * potential race with anon_decref. (segvn_unmap and
418             * segvn_free needs the as writer lock to do anon_free.)
419             */
420             if (*app != NULL) {
421                 #if 0
422                 if ((*app)->an_refcnt == 0)
423                 /*
424                 * Consider the following scenario (unlikey
425                 * though):
426                 * 1. an_refcnt == 2
427                 * 2. we softlock the page.
428                 * 3. cow occurs on this addr. So a new ap,
429                 * page and mapping is established on addr.
430                 * 4. an_refcnt drops to 1 (segvn_faultpage
431                 * -> anon_decref(oldap))
432                 * 5. the last ref to ap also drops (from
433                 * another as). It ends up blocked inside
434                 * anon_decref trying to get page's excl lock.
435                 * 6. Later kcfree unlocks the page, call
436                 * anon_decref -> oops, ap is gone already.
437                 *
438                 * Holding as writer lock solves all problems.
439                 */
440                 *app = NULL;
441             }
442             else

```

```

442 #endif
443                                     (*app)->an_refcnt++;
444                                     }
445                                     mutex_exit(ahm);
446     } else {
447         *app = NULL;
448     }
449     if (kaddr != (caddr_t)-1) {
450         if (pp != *cached_ppp) {
451             if (*cached_ppp == NULL)
452                 flags = HAT_LOAD_LOCK | HAT_NOSYNC |
453                       HAT_LOAD_NOCONSIST;
454             else
455                 flags = HAT_LOAD_REMAP |
456                       HAT_LOAD_NOCONSIST;
457             /*
458              * In order to cache the kernel mapping after
459              * the user page is unlocked, we call
460              * hat_devload instead of hat_memload so
461              * that the kernel mapping we set up here is
462              * "invisible" to the rest of the world. This
463              * is not very pretty. But as long as the
464              * caller bears the responsibility of keeping
465              * cache consistency, we should be ok -
466              * HAT_NOCONSIST will get us a uncached
467              * mapping on VAC. hat_softlock will flush
468              * a VAC_WRITEBACK cache. Therefore the kaddr
469              * doesn't have to be of the same vcolor as
470              * uaddr.
471              * The alternative is - change hat_devload
472              * to get a cached mapping. Allocate a kaddr
473              * with the same vcolor as uaddr. Then
474              * hat_softlock won't need to flush the VAC.
475              */
476             hat_devload(kas.a_hat, kaddr, PAGE_SIZE,
477                       page_pptonum(pp), PROT_READ, flags);
478             *cached_ppp = pp;
479         }
480         kaddr += PAGE_SIZE;
481     }
482     cached_ppp++;
483     app++;
484     ++i;
485 }
486 if (cow) {
487     AS_LOCK_EXIT(as, &as->a_lock);
488 }
489 if (first && res == FC_NOMAP) {
490     /*
491      * If the address is not mapped yet, we call as_fault to
492      * fault the pages in. We could've fallen back to copy and
493      * let it fault in the pages. But for a mapped file, we
494      * normally reference each page only once. For zero-copy to
495      * be of any use, we'd better fall in the page now and try
496      * again.
497      */
498     first = 0;
499     size = size << PAGESHIFT;
500     uaddr += size;
501     total -= size;
502     size = total;
503     res = as_fault(as->a_hat, as, uaddr, size, F_INVAL, S_READ);
504     if (cow)
505         AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
506     goto tryagain;
507 }

```

```

508     switch (res) {
509     case FC_NOSUPPORT:
510         return (ENOTSUP);
511     case FC_PROT: /* Pretend we don't know about it. This will be */
512                 /* caught by the caller when uiomove fails. */
513     case FC_NOMAP:
514     case FC_OBJERR:
515     default:
516         return (0);
517     }
518 }

```

unchanged portion omitted


```

*****
38864 Tue Nov 24 09:34:47 2015
new/usr/src/uts/common/os/watchpoint.c
patch lower-case-segops
*****
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, Version 1.0 only
6  * (the "License"). You may not use this file except in compliance
7  * with the License.
8  *
9  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
10 * or http://www.opensolaris.org/os/licensing.
11 * See the License for the specific language governing permissions
12 * and limitations under the License.
13 *
14 * When distributing Covered Code, include this CDDL HEADER in each
15 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
16 * If applicable, add the following below this CDDL HEADER, with the
17 * fields enclosed by brackets "[]" replaced with your own identifying
18 * information: Portions Copyright [yyyy] [name of copyright owner]
19 *
20 * CDDL HEADER END
21 */
22 /*
23 * Copyright 2005 Sun Microsystems, Inc. All rights reserved.
24 * Use is subject to license terms.
25 */

27 #pragma ident      "%Z%M% %I%      %E% SMI"

27 #include <sys/types.h>
28 #include <sys/t_lock.h>
29 #include <sys/param.h>
30 #include <sys/cred.h>
31 #include <sys/debug.h>
32 #include <sys/inline.h>
33 #include <sys/kmem.h>
34 #include <sys/proc.h>
35 #include <sys/regset.h>
36 #include <sys/sysmacros.h>
37 #include <sys/system.h>
38 #include <sys/prsystem.h>
39 #include <sys/buf.h>
40 #include <sys/signal.h>
41 #include <sys/user.h>
42 #include <sys/cpuvar.h>

44 #include <sys/fault.h>
45 #include <sys/syscall.h>
46 #include <sys/procfs.h>
47 #include <sys/cmn_err.h>
48 #include <sys/stack.h>
49 #include <sys/watchpoint.h>
50 #include <sys/copyops.h>
51 #include <sys/schedctl.h>

53 #include <sys/mman.h>
54 #include <vm/as.h>
55 #include <vm/seg.h>

57 /*
58 * Copy ops vector for watchpoints.
59 */

```

```

60 static int      watch_copyin(const void *, void *, size_t);
61 static int      watch_xcopyin(const void *, void *, size_t);
62 static int      watch_copyout(const void *, void *, size_t);
63 static int      watch_xcopyout(const void *, void *, size_t);
64 static int      watch_copyinstr(const char *, char *, size_t, size_t *);
65 static int      watch_copyoutstr(const char *, char *, size_t, size_t *);
66 static int      watch_fuword8(const void *, uint8_t *);
67 static int      watch_fuword16(const void *, uint16_t *);
68 static int      watch_fuword32(const void *, uint32_t *);
69 static int      watch_suword8(void *, uint8_t);
70 static int      watch_suword16(void *, uint16_t);
71 static int      watch_suword32(void *, uint32_t);
72 static int      watch_physio(int (*)(struct buf *), struct buf *,
73                          dev_t, int, void (*)(struct buf *), struct uio *);
74 #ifdef _LP64
75 static int      watch_fuword64(const void *, uint64_t *);
76 static int      watch_suword64(void *, uint64_t);
77 #endif

79 struct copyops watch_copyops = {
80     watch_copyin,
81     watch_xcopyin,
82     watch_copyout,
83     watch_xcopyout,
84     watch_copyinstr,
85     watch_copyoutstr,
86     watch_fuword8,
87     watch_fuword16,
88     watch_fuword32,
89 #ifdef _LP64
90     watch_fuword64,
91 #else
92     NULL,
93 #endif
94     watch_suword8,
95     watch_suword16,
96     watch_suword32,
97 #ifdef _LP64
98     watch_suword64,
99 #else
100    NULL,
101 #endif
102     watch_physio
103 };
_____unchanged_portion_omitted_____

151 #define X      0
152 #define W      1
153 #define R      2
154 #define sum(a) (a[X] + a[W] + a[R])

156 /*
157 * Common code for pr_mappage() and pr_unmappage().
158 */
159 static int
160 pr_do_mappage(caddr_t addr, size_t size, int mapin, enum seg_rw rw, int kernel)
161 {
162     proc_t *p = curproc;
163     struct as *as = p->p_as;
164     char *eaddr = addr + size;
165     int prot_rw = rw_to_prot(rw);
166     int xrw = rw_to_index(rw);
167     int rv = 0;
168     struct watched_page *pwp;
169     struct watched_page tpw;
170     avl_index_t where;

```

```

171     uint_t prot;
173     ASSERT(as != &kas);

175 startover:
176     ASSERT(rv == 0);
177     if (avl_numnodes(&as->a_wpage) == 0)
178         return (0);

180 /*
181  * as->a_wpage can only be changed while the process is totally stopped.
182  * Don't grab p_lock here. Holding p_lock while grabbing the address
183  * space lock leads to deadlocks with the clock thread. Note that if an
184  * as_fault() is servicing a fault to a watched page on behalf of an
185  * XHAT provider, watchpoint will be temporarily cleared (and wp_prot
186  * will be set to wp_oprot). Since this is done while holding as writer
187  * lock, we need to grab as lock (reader lock is good enough).
188  *
189  * p_maplock prevents simultaneous execution of this function. Under
190  * normal circumstances, holdwatch() will stop all other threads, so the
191  * lock isn't really needed. But there may be multiple threads within
192  * stop() when SWATCHOK is set, so we need to handle multiple threads
193  * at once. See holdwatch() for the details of this dance.
194  */

196     mutex_enter(&p->p_maplock);
197     AS_LOCK_ENTER(as, &as->a_lock, RW_READER);

199     tpw.wp_vaddr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);
200     if ((pwp = avl_find(&as->a_wpage, &tpw, &where)) == NULL)
201         pwp = avl_nearest(&as->a_wpage, where, AVL_AFTER);

203     for (; pwp != NULL && pwp->wp_vaddr < eaddr;
204          pwp = AVL_NEXT(&as->a_wpage, pwp)) {
206         /*
207          * If the requested protection has not been
208          * removed, we need not remap this page.
209          */
210         prot = pwp->wp_prot;
211         if (kernel || (prot & PROT_USER))
212             if (prot & prot_rw)
213                 continue;
214         /*
215          * If the requested access does not exist in the page's
216          * original protections, we need not remap this page.
217          * If the page does not exist yet, we can't test it.
218          */
219         if ((prot = pwp->wp_oprot) != 0) {
220             if (!(kernel || (prot & PROT_USER)))
221                 continue;
222             if (!(prot & prot_rw))
223                 continue;
224         }

226         if (mapin) {
227             /*
228              * Before mapping the page in, ensure that
229              * all other lwps are held in the kernel.
230              */
231             if (p->p_mapcnt == 0) {
232                 /*
233                  * Release as lock while in holdwatch()
234                  * in case other threads need to grab it.
235                  */
236                 AS_LOCK_EXIT(as, &as->a_lock);

```

```

237         mutex_exit(&p->p_maplock);
238         if (holdwatch() != 0) {
239             /*
240              * We stopped in holdwatch().
241              * Start all over again because the
242              * watched page list may have changed.
243              */
244             goto startover;
245         }
246         mutex_enter(&p->p_maplock);
247         AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
248     }
249     p->p_mapcnt++;
250 }

252     addr = pwp->wp_vaddr;
253     rv++;

255     prot = pwp->wp_prot;
256     if (mapin) {
257         if (kernel)
258             pwp->wp_kmap[xrw]++;
259         else
260             pwp->wp_umap[xrw]++;
261         pwp->wp_flags |= WP_NOWATCH;
262         if (pwp->wp_kmap[X] + pwp->wp_umap[X])
263             /* cannot have exec-only protection */
264             prot |= PROT_READ|PROT_EXEC;
265         if (pwp->wp_kmap[R] + pwp->wp_umap[R])
266             prot |= PROT_READ;
267         if (pwp->wp_kmap[W] + pwp->wp_umap[W])
268             /* cannot have write-only protection */
269             prot |= PROT_READ|PROT_WRITE;
270     #if 0 /* damned broken mmu feature! */
271         if (sum(pwp->wp_umap) == 0)
272             prot &= ~PROT_USER;
273     #endif
274     } else {
275         ASSERT(pwp->wp_flags & WP_NOWATCH);
276         if (kernel) {
277             ASSERT(pwp->wp_kmap[xrw] != 0);
278             --pwp->wp_kmap[xrw];
279         } else {
280             ASSERT(pwp->wp_umap[xrw] != 0);
281             --pwp->wp_umap[xrw];
282         }
283         if (sum(pwp->wp_kmap) + sum(pwp->wp_umap) == 0)
284             pwp->wp_flags &= ~WP_NOWATCH;
285         else {
286             if (pwp->wp_kmap[X] + pwp->wp_umap[X])
287                 /* cannot have exec-only protection */
288                 prot |= PROT_READ|PROT_EXEC;
289             if (pwp->wp_kmap[R] + pwp->wp_umap[R])
290                 prot |= PROT_READ;
291             if (pwp->wp_kmap[W] + pwp->wp_umap[W])
292                 /* cannot have write-only protection */
293                 prot |= PROT_READ|PROT_WRITE;
294     #if 0 /* damned broken mmu feature! */
295             if (sum(pwp->wp_umap) == 0)
296                 prot &= ~PROT_USER;
297     #endif
298         }
299     }

302     if (pwp->wp_oprot != 0) { /* if page exists */

```

```

303         struct seg *seg;
304         uint_t oprot;
305         int err, retrycnt = 0;

307         AS_LOCK_EXIT(as, &as->a_lock);
308         AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
309     retry:
310         seg = as_segat(as, addr);
311         ASSERT(seg != NULL);
312         (void) segop_getprot(seg, addr, 0, &oprot);
313         SEGOP_GETPROT(seg, addr, 0, &oprot);
314         if (prot != oprot) {
315             err = segop_setprot(seg, addr, PAGESIZE, prot);
316             err = SEGOP_SETPROT(seg, addr, PAGESIZE, prot);
317             if (err == IE_RETRY) {
318                 ASSERT(retrycnt == 0);
319                 retrycnt++;
320                 goto retry;
321             }
322         }
323     } else
324         AS_LOCK_EXIT(as, &as->a_lock);

325     /*
326     * When all pages are mapped back to their normal state,
327     * continue the other lwps.
328     */
329     if (!mapin) {
330         ASSERT(p->p_mapcnt > 0);
331         p->p_mapcnt--;
332         if (p->p_mapcnt == 0) {
333             mutex_exit(&p->p_maplock);
334             mutex_enter(&p->p_lock);
335             continuelwps(p);
336             mutex_exit(&p->p_lock);
337             mutex_enter(&p->p_maplock);
338         }
339     }

341     AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
342 }

344     AS_LOCK_EXIT(as, &as->a_lock);
345     mutex_exit(&p->p_maplock);

347     return (rv);
348 }

    unchanged_portion_omitted

376 /*
377  * Function called by an lwp after it resumes from stop().
378  */
379 void
380 setallwatch(void)
381 {
382     proc_t *p = curproc;
383     struct as *as = curproc->p_as;
384     struct watched_page *pwp, *next;
385     struct seg *seg;
386     caddr_t vaddr;
387     uint_t prot;
388     int err, retrycnt;

390     if (p->p_wprot == NULL)
391         return;

```

```

393     ASSERT(MUTEX_NOT_HELD(&curproc->p_lock));

395     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);

397     pwp = p->p_wprot;
398     while (pwp != NULL) {

400         vaddr = pwp->wp_vaddr;
401         retrycnt = 0;
402     retry:
403         ASSERT(pwp->wp_flags & WP_SETPROT);
404         if ((seg = as_segat(as, vaddr)) != NULL &&
405             !(pwp->wp_flags & WP_NOWATCH)) {
406             prot = pwp->wp_prot;
407             err = segop_setprot(seg, vaddr, PAGESIZE, prot);
408             err = SEGOP_SETPROT(seg, vaddr, PAGESIZE, prot);
409             if (err == IE_RETRY) {
410                 ASSERT(retrycnt == 0);
411                 retrycnt++;
412                 goto retry;
413             }
414         }

415         next = pwp->wp_list;

417         if (pwp->wp_read + pwp->wp_write + pwp->wp_exec == 0) {
418             /*
419             * No watched areas remain in this page.
420             * Free the watched_page structure.
421             */
422             avl_remove(&as->a_wpage, pwp);
423             kmem_free(pwp, sizeof (struct watched_page));
424         } else {
425             pwp->wp_flags &= ~WP_SETPROT;
426         }

428         pwp = next;
429     }
430     p->p_wprot = NULL;

432     AS_LOCK_EXIT(as, &as->a_lock);
433 }

    unchanged_portion_omitted

```

```
*****
193225 Tue Nov 24 09:34:47 2015
new/usr/src/uts/common/os/zone.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

5593 /*
5594  * Return zero if the process has at least one vnode mapped in to its
5595  * address space which shouldn't be allowed to change zones.
5596  *
5597  * Also return zero if the process has any shared mappings which reserve
5598  * swap. This is because the counting for zone.max-swap does not allow swap
5599  * reservation to be shared between zones. zone swap reservation is counted
5600  * on zone->zone_max_swap.
5601  */
5602 static int
5603 as_can_change_zones(void)
5604 {
5605     proc_t *pp = curproc;
5606     struct seg *seg;
5607     struct as *as = pp->p_as;
5608     vnode_t *vp;
5609     int allow = 1;

5611     ASSERT(pp->p_as != &kas);
5612     AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
5613     for (seg = AS_SEGFIRST(as); seg != NULL; seg = AS_SEGNEXT(as, seg)) {

5615         /*
5616          * Cannot enter zone with shared anon memory which
5617          * reserves swap. See comment above.
5618          */
5619         if (seg_can_change_zones(seg) == B_FALSE) {
5620             allow = 0;
5621             break;
5622         }
5623         /*
5624          * if we can't get a backing vnode for this segment then skip
5625          * it.
5626          */
5627         vp = NULL;
5628         if (segop_getvp(seg, seg->s_base, &vp) != 0 || vp == NULL)
5628         if (SEGOP_GETVP(seg, seg->s_base, &vp) != 0 || vp == NULL)
5629             continue;
5630         if (!vn_can_change_zones(vp)) { /* bail on first match */
5631             allow = 0;
5632             break;
5633         }
5634     }
5635     AS_LOCK_EXIT(as, &as->a_lock);
5636     return (allow);
5637 }
_____unchanged_portion_omitted_____
```

new/usr/src/uts/common/sys/watchpoint.h

1

```
*****
4198 Tue Nov 24 09:34:47 2015
new/usr/src/uts/common/sys/watchpoint.h
patch lower-case-segops
*****
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, Version 1.0 only
6  * (the "License"). You may not use this file except in compliance
7  * with the License.
8  *
9  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
10 * or http://www.opensolaris.org/os/licensing.
11 * See the License for the specific language governing permissions
12 * and limitations under the License.
13 *
14 * When distributing Covered Code, include this CDDL HEADER in each
15 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
16 * If applicable, add the following below this CDDL HEADER, with the
17 * fields enclosed by brackets "[]" replaced with your own identifying
18 * information: Portions Copyright [yyyy] [name of copyright owner]
19 *
20 * CDDL HEADER END
21 */
22 /*
23  * Copyright 2004 Sun Microsystems, Inc. All rights reserved.
24  * Use is subject to license terms.
25  */

27 #ifndef _SYS_WATCHPOINT_H
28 #define _SYS_WATCHPOINT_H

30 #pragma ident      "%Z%%M% %I%      %E% SMI"

30 #include <sys/types.h>
31 #include <vm/seg_enum.h>
32 #include <sys/copyops.h>
33 #include <sys/avl.h>

35 #ifdef __cplusplus
36 extern "C" {
37 #endif

39 /*
40  * Definitions for the VM implementation of watchpoints.
41  * See proc(4) and <sys/procfs.h> for definitions of the user interface.
42  */

44 /*
45  * Each process with watchpoints has a linked list of watched areas.
46  * The list is kept sorted by user-level virtual address.
47  */
48 typedef struct watched_area {
49     avl_node_t wa_link;      /* link in AVL tree */
50     caddr_t wa_vaddr;       /* virtual address of watched area */
51     caddr_t wa_eaddr;       /* virtual address plus size */
52     ulong_t wa_flags;       /* watch type flags (see <sys/procfs.h>) */
53 } watched_area_t;
unchanged_portion_omitted

73 /* wp_flags */
74 #define WP_NOWATCH      0x01 /* protections temporarily restored */
75 #define WP_SETPROT      0x02 /* segop_setprot() needed on this page */
77 #define WP_SETPROT      0x02 /* SEGOP_SETPROT() needed on this page */
```

new/usr/src/uts/common/sys/watchpoint.h

2

```
77 #ifdef _KERNEL

79 /*
80  * These functions handle the necessary logic to perform the copy operation
81  * while ignoring watchpoints.
82  */
83 extern int copyin_nowatch(const void *, void *, size_t);
84 extern int copyout_nowatch(const void *, void *, size_t);
85 extern int fuword32_nowatch(const void *, uint32_t *);
86 extern int suword32_nowatch(void *, uint32_t);
87 #ifdef _LP64
88 extern int suword64_nowatch(void *, uint64_t);
89 extern int fuword64_nowatch(const void *, uint64_t *);
90 #endif

92 /*
93  * Disable watchpoints for a given region of memory. When bracketed by these
94  * calls, functions can use copyops and ignore watchpoints.
95  */
96 extern int watch_disable_addr(const void *, size_t, enum seg_rw);
97 extern void watch_enable_addr(const void *, size_t, enum seg_rw);

99 /*
100  * Enable/Disable watchpoints for an entire thread.
101  */
102 extern void watch_enable(kthread_id_t);
103 extern void watch_disable(kthread_id_t);

105 struct as;
106 struct proc;
107 struct k_siginfo;
108 extern void setallwatch(void);
109 extern int pr_is_watchpage(caddr_t, enum seg_rw);
110 extern int pr_is_watchpage_as(caddr_t, enum seg_rw, struct as *);
111 extern int pr_is_watchpoint(caddr_t *, int *, size_t, size_t *,
112                             enum seg_rw);
113 extern void do_watch_step(caddr_t, size_t, enum seg_rw, int, greg_t);
114 extern int undo_watch_step(struct k_siginfo *);
115 extern int wp_compare(const void *, const void *);
116 extern int wa_compare(const void *, const void *);

118 extern struct copyops watch_copyops;

120 extern watched_area_t *pr_find_watched_area(struct proc *, watched_area_t *,
121                                             avl_index_t *);

123 #endif

125 #ifdef __cplusplus
126 }
unchanged_portion_omitted
```

new/usr/src/uts/common/syscall/utssys.c

1

```
*****
23713 Tue Nov 24 09:34:47 2015
new/usr/src/uts/common/syscall/utssys.c
patch lower-case-segops
*****
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, Version 1.0 only
6  * (the "License"). You may not use this file except in compliance
7  * with the License.
8  *
9  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
10 * or http://www.opensolaris.org/os/licensing.
11 * See the License for the specific language governing permissions
12 * and limitations under the License.
13 *
14 * When distributing Covered Code, include this CDDL HEADER in each
15 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
16 * If applicable, add the following below this CDDL HEADER, with the
17 * fields enclosed by brackets "[]" replaced with your own identifying
18 * information: Portions Copyright [yyyy] [name of copyright owner]
19 *
20 * CDDL HEADER END
21 */

23 /*
24  * Copyright 2004 Sun Microsystems, Inc. All rights reserved.
25  * Use is subject to license terms.
26  */

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

32 #pragma ident      "%Z%M% %I%      %E% SMI"

32 #include <sys/param.h>
33 #include <sys/inttypes.h>
34 #include <sys/types.h>
35 #include <sys/sysmacros.h>
36 #include <sys/systm.h>
37 #include <sys/user.h>
38 #include <sys/errno.h>
39 #include <sys/vfs.h>
40 #include <sys/vnode.h>
41 #include <sys/file.h>
42 #include <sys/proc.h>
43 #include <sys/session.h>
44 #include <sys/var.h>
45 #include <sys/utsname.h>
46 #include <sys/utssys.h>
47 #include <sys/ustat.h>
48 #include <sys/statvfs.h>
49 #include <sys/kmem.h>
50 #include <sys/debug.h>
51 #include <sys/pathname.h>
52 #include <sys/modctl.h>
53 #include <sys/fs/snode.h>
54 #include <sys/sunldi_impl.h>
55 #include <sys/ddi.h>
56 #include <sys/sunddi.h>
57 #include <sys/cmn_err.h>
58 #include <sys/ddipropdefs.h>
59 #include <sys/ddi_impldefs.h>
```

new/usr/src/uts/common/syscall/utssys.c

2

```
60 #include <sys/modctl.h>
61 #include <sys/flock.h>
62 #include <sys/share.h>
63 #include <vm/as.h>
64 #include <vm/seg.h>
65 #include <vm/seg_vn.h>
66 #include <util/qsorth.h>
67 #include <sys/zone.h>

69 /*
70  * utssys()
71  */
72 static int      uts_fusers(char *, int, intptr_t);
73 static int      _statvfs64_by_dev(dev_t, struct statvfs64 *);

75 #if defined(_ILP32) || defined(_SYSCALL32_IMPL)

77 static int utssys_uname32(caddr_t, rval_t *);
78 static int utssys_ustat32(dev_t, struct ustat32 *);

80 int64_t
81 utssys32(void *buf, int arg, int type, void *outbp)
82 {
83     int error;
84     rval_t rv;

86     rv.r_vals = 0;

88     switch (type) {
89     case UTS_UNAME:
90         /*
91          * This is an obsolete way to get the utsname structure
92          * (it only gives you the first 8 characters of each field!)
93          * uname(2) is the preferred and better interface.
94          */
95         error = utssys_uname32(buf, &rv);
96         break;
97     case UTS_USTAT:
98         error = utssys_ustat32(expldev((dev32_t)arg), buf);
99         break;
100    case UTS_FUSERS:
101        error = utssys_fusers(buf, arg, (intptr_t)outbp);
102        break;
103    default:
104        error = EINVAL;
105        break;
106    }

108    return (error == 0 ? rv.r_vals : (int64_t)set_errno(error));
109 }

_____ unchanged_portion_omitted _____

311 static fu_data_t *
312 dofusers(vnode_t *fvp, int flags)
313 {
314     fu_data_t      *fu_data;
315     proc_t         *prp;
316     vfs_t          *cvfsp;
317     pid_t          npids, pidx, *pidlist;
318     int            v_proc = v.v_proc;      /* max # of procs */
319     int            pcnt = 0;
320     int            contained = (flags & F_CONTAINED);
321     int            nbmandonly = (flags & F_NBMANDLIST);
322     int            dip_usage = (flags & F_DEVINFO);
323     int            fvp_isdev = vn_matchchops(fvp, spec_getvnnodeops());
324     zone_t *zone = curproc->p_zone;
```

```

325     int inglobal = INGLOBALZONE(curproc);

327     /* get a pointer to the file system containing this vnode */
328     cvfsp = fvp->v_vfsp;
329     ASSERT(cvfsp);

331     /* allocate the data structure to return our results in */
332     fu_data = kmem_alloc(fu_data_size(v_proc), KM_SLEEP);
333     fu_data->fud_user_max = v_proc;
334     fu_data->fud_user_count = 0;

336     /* get a snapshot of all the pids we're going to check out */
337     pidlist = kmem_alloc(v_proc * sizeof(pid_t), KM_SLEEP);
338     mutex_enter(&pidlock);
339     for (npids = 0, prp = practive; prp != NULL; prp = prp->p_next) {
340         if (inglobal || prp->p_zone == zone)
341             pidlist[npids++] = prp->p_pid;
342     }
343     mutex_exit(&pidlock);

345     /* grab each process and check its file usage */
346     for (pidx = 0; pidx < npids; pidx++) {
347         locklist_t      *llp = NULL;
348         uf_info_t        *fip;
349         vnode_t          *vp;
350         user_t           *up;
351         sess_t           *sp;
352         uid_t            uid;
353         pid_t            pid = pidlist[pidx];
354         int              i, use_flag = 0;

356         /*
357          * grab prp->p_lock using sprlock()
358          * if sprlock() fails the process does not exist anymore
359          */
360         prp = sprlock(pid);
361         if (prp == NULL)
362             continue;

364         /* get the processes credential info in case we need it */
365         mutex_enter(&prp->p_crlock);
366         uid = crgetruid(prp->p_cred);
367         mutex_exit(&prp->p_crlock);

369         /*
370          * it's safe to drop p_lock here because we
371          * called sprlock() before and it set the SPRLOCK
372          * flag for the process so it won't go away.
373          */
374         mutex_exit(&prp->p_lock);

376         /*
377          * now we want to walk a processes open file descriptors
378          * to do this we need to grab the fip->fi_lock. (you
379          * can't hold p_lock when grabbing the fip->fi_lock.)
380          */
381         fip = P_FINFO(prp);
382         mutex_enter(&fip->fi_lock);

384         /*
385          * Snapshot nbmand locks for pid
386          */
387         llp = flk_active_nbmand_locks(prp->p_pid);
388         for (i = 0; i < fip->fi_nfiles; i++) {
389             uf_entry_t    *ufp;
390             file_t        *fp;

```

```

392         UF_ENTER(ufp, fip, i);
393         if (((fip = ufp->uf_file) == NULL) ||
394             ((vp = fp->f_vnode) == NULL)) {
395             UF_EXIT(ufp);
396             continue;
397         }

399         /*
400          * if the target file (fvp) is not a device
401          * and corresponds to the root of a filesystem
402          * (cvfsp), then check if it contains the file
403          * is use by this process (vp).
404          */
405         if (contained && (vp->v_vfsp == cvfsp))
406             use_flag |= F_OPEN;

408         /*
409          * if the target file (fvp) is not a device,
410          * then check if it matches the file in use
411          * by this process (vp).
412          */
413         if (!fvp_isdev && VN_CMP(fvp, vp))
414             use_flag |= F_OPEN;

416         /*
417          * if the target file (fvp) is a device,
418          * then check if the current file in use
419          * by this process (vp) maps to the same device
420          * minor node.
421          */
422         if (fvp_isdev &&
423             vn_matchops(vp, spec_getvnodeops()) &&
424             (fvp->v_rdev == vp->v_rdev))
425             use_flag |= F_OPEN;

427         /*
428          * if the target file (fvp) is a device,
429          * and we're checking for device instance
430          * usage, then check if the current file in use
431          * by this process (vp) maps to the same device
432          * instance.
433          */
434         if (dip_usage &&
435             vn_matchops(vp, spec_getvnodeops()) &&
436             (VTOCS(fvp)->s_dip == VTOCS(vp)->s_dip))
437             use_flag |= F_OPEN;

439         /*
440          * if the current file in use by this process (vp)
441          * doesn't match what we're looking for, move on
442          * to the next file in the process.
443          */
444         if ((use_flag & F_OPEN) == 0) {
445             UF_EXIT(ufp);
446             continue;
447         }

449         if (proc_has_nbmand_on_vp(vp, prp->p_pid, llp)) {
450             /* A nbmand found so we're done. */
451             use_flag |= F_NBM;
452             UF_EXIT(ufp);
453             break;
454         }
455         UF_EXIT(ufp);
456     }

```

```

457     if (llp)
458         flk_free_locklist(llp);
460     mutex_exit(&fip->fi_lock);
462     /*
463     * If nbmand usage tracking is desired and no nbmand was
464     * found for this process, then no need to do further
465     * usage tracking for this process.
466     */
467     if (nbmandonly && (!(use_flag & F_NBM))) {
468         /*
469         * grab the process lock again, clear the SPRLOCK
470         * flag, release the process, and continue.
471         */
472         mutex_enter(&prp->p_lock);
473         sprunlock(prp);
474         continue;
475     }
477     /*
478     * All other types of usage.
479     * For the next few checks we need to hold p_lock.
480     */
481     mutex_enter(&prp->p_lock);
482     up = PTOU(prp);
483     if (fvp_isdev) {
484         /*
485         * if the target file (fvp) is a device
486         * then check if it matches the processes tty
487         *
488         * we grab s_lock to protect ourselves against
489         * freectty() freeing the vnode out from under us.
490         */
491         sp = prp->p_sessp;
492         mutex_enter(&sp->s_lock);
493         vp = prp->p_sessp->s_vp;
494         if (vp != NULL) {
495             if (fvp->v_rdev == vp->v_rdev)
496                 use_flag |= F_TTY;
498             if (dip_usage &&
499                 (VTOCS(fvp)->s_dip == VTOCS(vp)->s_dip))
500                 use_flag |= F_TTY;
501         }
502         mutex_exit(&sp->s_lock);
503     } else {
504         /* check the processes current working directory */
505         if (up->u_cdir &&
506             (VN_CMP(fvp, up->u_cdir) ||
507              (contained && (up->u_cdir->v_vfsp == cvfsp))))
508             use_flag |= F_CDIRE;
510         /* check the processes root directory */
511         if (up->u_rdir &&
512             (VN_CMP(fvp, up->u_rdir) ||
513              (contained && (up->u_rdir->v_vfsp == cvfsp))))
514             use_flag |= F_RDIRE;
516         /* check the program text vnode */
517         if (prp->p_exec &&
518             (VN_CMP(fvp, prp->p_exec) ||
519              (contained && (prp->p_exec->v_vfsp == cvfsp))))
520             use_flag |= F_TEXT;
521     }

```

```

523     /* Now we can drop p_lock again */
524     mutex_exit(&prp->p_lock);
526     /*
527     * now we want to walk a processes memory mappings.
528     * to do this we need to grab the prp->p_as lock. (you
529     * can't hold p_lock when grabbing the prp->p_as lock.)
530     */
531     if (prp->p_as != &kas) {
532         struct seg *seg;
533         struct as *as = prp->p_as;
535         AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
536         for (seg = AS_SEGFIRST(as); seg;
537              seg = AS_SEGNEXT(as, seg)) {
538             /*
539             * if we can't get a backing vnode for this
540             * segment then skip it
541             */
542             vp = NULL;
543             if ((segop_getvp(seg, seg->s_base, &vp)) ||
544                 if ((SEGOP_GETVP(seg, seg->s_base, &vp)) ||
545                     (vp == NULL))
546                 continue;
547             /*
548             * if the target file (fvp) is not a device
549             * and corresponds to the root of a filesystem
550             * (cvfsp), then check if it contains the
551             * vnode backing this segment (vp).
552             */
553             if (contained && (vp->v_vfsp == cvfsp)) {
554                 use_flag |= F_MAP;
555                 break;
556             }
558             /*
559             * if the target file (fvp) is not a device,
560             * check if it matches the the vnode backing
561             * this segment (vp).
562             */
563             if (!fvp_isdev && VN_CMP(fvp, vp)) {
564                 use_flag |= F_MAP;
565                 break;
566             }
568             /*
569             * if the target file (fvp) isn't a device,
570             * or the the vnode backing this segment (vp)
571             * isn't a device then continue.
572             */
573             if (!fvp_isdev ||
574                 !vn_matchops(vp, spec_getvnodeops()))
575                 continue;
577             /*
578             * check if the vnode backing this segment
579             * (vp) maps to the same device minor node
580             * as the target device (fvp)
581             */
582             if (fvp->v_rdev == vp->v_rdev) {
583                 use_flag |= F_MAP;
584                 break;
585             }
587             /*

```



```
588         * if we're checking for device instance
589         * usage, then check if the vnode backing
590         * this segment (vp) maps to the same device
591         * instance as the target device (fvp).
592         */
593         if (dip_usage &&
594             (VTOCS(fvp)->s_dip == VTOCS(vp)->s_dip)) {
595             use_flag |= F_MAP;
596             break;
597         }
598     }
599     AS_LOCK_EXIT(as, &as->a_lock);
600 }
601
602 if (use_flag) {
603     ASSERT(pcnt < fu_data->fud_user_max);
604     fu_data->fud_user[pcnt].fu_flags = use_flag;
605     fu_data->fud_user[pcnt].fu_pid = pid;
606     fu_data->fud_user[pcnt].fu_uid = uid;
607     pcnt++;
608 }
609
610 /*
611  * grab the process lock again, clear the SPRLOCK
612  * flag, release the process, and continue.
613  */
614 mutex_enter(&prp->p_lock);
615 sprunlock(prp);
616 }
617
618 kmem_free(pidlist, v_proc * sizeof (pid_t));
619
620 fu_data->fud_user_count = pcnt;
621 return (fu_data);
622 }
unchanged_portion_omitted
```

new/usr/src/uts/common/vm/seg.h

1

```
*****
10256 Tue Nov 24 09:34:48 2015
new/usr/src/uts/common/vm/seg.h
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

147 #ifndef _KERNEL

149 /*
150  * Generic segment operations
151  */
152 extern void seg_init(void);
153 extern struct seg *seg_alloc(struct as *as, caddr_t base, size_t size);
154 extern int seg_attach(struct as *as, caddr_t base, size_t size,
155 struct seg *seg);
156 extern void seg_unmap(struct seg *seg);
157 extern void seg_free(struct seg *seg);

159 /*
160  * functions for pagelock cache support
161  */
162 typedef int (*seg_preclaim_cbfunc_t)(void *, caddr_t, size_t,
163 struct page **, enum seg_rw, int);

165 extern struct page **seg_plookup(struct seg *seg, struct anon_map *amp,
166 caddr_t addr, size_t len, enum seg_rw rw, uint_t flags);
167 extern void seg_pinactive(struct seg *seg, struct anon_map *amp,
168 caddr_t addr, size_t len, struct page **pp, enum seg_rw rw,
169 uint_t flags, seg_preclaim_cbfunc_t callback);

171 extern void seg_ppurge(struct seg *seg, struct anon_map *amp,
172 uint_t flags);
173 extern void seg_ppurge_wiredpp(struct page **pp);

175 extern int seg_pininsert_check(struct seg *seg, struct anon_map *amp,
176 caddr_t addr, size_t len, uint_t flags);
177 extern int seg_pininsert(struct seg *seg, struct anon_map *amp,
178 caddr_t addr, size_t len, size_t wlen, struct page **pp, enum seg_rw rw,
179 uint_t flags, seg_preclaim_cbfunc_t callback);

181 extern void seg_pasync_thread(void);
182 extern void seg_preap(void);
183 extern int seg_p_disable(void);
184 extern void seg_p_enable(void);

186 extern segadvstat_t segadvstat;

188 /*
189  * Flags for pagelock cache support.
190  * Flags argument is passed as uint_t to pcache routines. upper 16 bits of
191  * the flags argument are reserved for alignment page shift when SEGP_PSHIFT
192  * is set.
193  */
194 #define SEGP_FORCE_WIRED 0x1 /* skip check against seg_pwindow */
195 #define SEGP_AMP 0x2 /* anon map's pcache entry */
196 #define SEGP_PSHIFT 0x4 /* addr pgsz shift for hash function */

198 /*
199  * Return values for seg_pininsert and seg_pininsert_check functions.
200  */
201 #define SEGP_SUCCESS 0 /* seg_pininsert() succeeded */
202 #define SEGP_FAIL 1 /* seg_pininsert() failed */

204 /* Page status bits for segop_incore */
205 #define SEG_PAGE_INCORE 0x01 /* VA has a page backing it */
```

new/usr/src/uts/common/vm/seg.h

2

```
206 #define SEG_PAGE_LOCKED 0x02 /* VA has a page that is locked */
207 #define SEG_PAGE_HASCOW 0x04 /* VA has a page with a copy-on-write */
208 #define SEG_PAGE_SOFTLOCK 0x08 /* VA has a page with softlock held */
209 #define SEG_PAGE_VNODEBACKED 0x10 /* Segment is backed by a vnode */
210 #define SEG_PAGE_ANON 0x20 /* VA has an anonymous page */
211 #define SEG_PAGE_VNODE 0x40 /* VA has a vnode page backing it */

213 #define seg_page(seg, addr) \
214 ((uintptr_t)((addr) - (seg)->s_base)) >> PAGESHIFT

216 #define seg_pages(seg) \
217 ((uintptr_t)((seg)->s_size + PAGEOFFSET) >> PAGESHIFT)

219 #define IE_NOMEM -1 /* internal to seg layer */
220 #define IE_RETRY -2 /* internal to seg layer */
221 #define IE_REATTACH -3 /* internal to seg layer */

223 /* Values for segop_inherit */
223 /* Values for SEGOP_INHERIT */
224 #define SEGP_INH_ZERO 0x01

226 int seg_inherit_notsup(struct seg *, caddr_t, size_t, uint_t);

228 /* Delay/retry factors for seg_p_mem_config_pre_del */
229 #define SEGP_PREDEL_DELAY_FACTOR 4
230 /*
231  * As a workaround to being unable to purge the pagelock
232  * cache during a DR delete memory operation, we use
233  * a stall threshold that is twice the maximum seen
234  * during testing. This workaround will be removed
235  * when a suitable fix is found.
236  */
237 #define SEGP_STALL_SECONDS 25
238 #define SEGP_STALL_THRESHOLD \
239 (SEGP_STALL_SECONDS * SEGP_PREDEL_DELAY_FACTOR)

241 #ifndef VMDEBUG

243 uint_t seg_page(struct seg *, caddr_t);
244 uint_t seg_pages(struct seg *);

246 #endif /* VMDEBUG */

248 boolean_t seg_can_change_zones(struct seg *);
249 size_t seg_swresv(struct seg *);

251 /* segop wrappers */
252 extern int segop_dup(struct seg *, struct seg *);
253 extern int segop_unmap(struct seg *, caddr_t, size_t);
254 extern void segop_free(struct seg *);
255 extern faultcode_t segop_fault(struct hat *, struct seg *, caddr_t, size_t,
256 enum fault_type, enum seg_rw);
257 extern faultcode_t segop_faulta(struct seg *, caddr_t);
258 extern int segop_setprot(struct seg *, caddr_t, size_t, uint_t);
259 extern int segop_checkprot(struct seg *, caddr_t, size_t, uint_t);
260 extern int segop_kluster(struct seg *, caddr_t, ssize_t);
261 extern size_t segop_swapout(struct seg *);
262 extern int segop_sync(struct seg *, caddr_t, size_t, int, uint_t);
263 extern size_t segop_incore(struct seg *, caddr_t, size_t, char *);
264 extern int segop_lockop(struct seg *, caddr_t, size_t, int, int, ulong_t *,
265 size_t);
266 extern int segop_getprot(struct seg *, caddr_t, size_t, uint_t *);
267 extern u_offset_t segop_getoffset(struct seg *, caddr_t);
268 extern int segop_gettype(struct seg *, caddr_t);
269 extern int segop_getvp(struct seg *, caddr_t, struct vnode **);
270 extern int segop_advise(struct seg *, caddr_t, size_t, uint_t);
```

```
271 extern void segop_dump(struct seg *);
272 extern int segop_pagelock(struct seg *, caddr_t, size_t, struct page ***,
273     enum lock_type, enum seg_rw);
274 extern int segop_setpagesize(struct seg *, caddr_t, size_t, uint_t);
275 extern int segop_getmemid(struct seg *, caddr_t, memid_t *);
276 extern struct lgrp_mem_policy_info *segop_getpolicy(struct seg *, caddr_t);
277 extern int segop_capable(struct seg *, segcapability_t);
278 extern int segop_inherit(struct seg *, caddr_t, size_t, uint_t);

280 #endif /* _KERNEL */

282 #ifdef __cplusplus
283 }
unchanged_portion_omitted
```

```

*****
114107 Tue Nov 24 09:34:48 2015
new/usr/src/uts/common/vm/seg_dev.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

355 /*
356  * Create a device segment.
357  */
358 int
359 segdev_create(struct seg *seg, void *argsp)
360 {
361     struct segdev_data *sdp;
362     struct segdev_crargs *a = (struct segdev_crargs *)argsp;
363     devmap_handle_t *dhp = (devmap_handle_t *)a->devmap_data;
364     int error;

366     /*
367      * Since the address space is "write" locked, we
368      * don't need the segment lock to protect "segdev" data.
369      */
370     ASSERT(seg->s_as && AS_WRITE_HELD(seg->s_as, &seg->s_as->a_lock));

372     hat_map(seg->s_as->a_hat, seg->s_base, seg->s_size, HAT_MAP);

374     sdp = sdp_alloc();

376     sdp->mapfunc = a->mapfunc;
377     sdp->offset = a->offset;
378     sdp->prot = a->prot;
379     sdp->maxprot = a->maxprot;
380     sdp->type = a->type;
381     sdp->pageprot = 0;
382     sdp->softlockcnt = 0;
383     sdp->vpage = NULL;

385     if (sdp->mapfunc == NULL)
386         sdp->devmap_data = dhp;
387     else
388         sdp->devmap_data = dhp = NULL;

390     sdp->hat_flags = a->hat_flags;
391     sdp->hat_attr = a->hat_attr;

393     /*
394      * Currently, hat_flags supports only HAT_LOAD_NOCONSIST
395      */
396     ASSERT(!(sdp->hat_flags & ~HAT_LOAD_NOCONSIST));

398     /*
399      * Hold shadow vnode -- segdev only deals with
400      * character (VCHR) devices. We use the common
401      * vp to hang pages on.
402      */
403     sdp->vp = specfind(a->dev, VCHR);
404     ASSERT(sdp->vp != NULL);

406     seg->s_ops = &segdev_ops;
407     seg->s_data = sdp;

409     while (dhp != NULL) {
410         dhp->dh_seg = seg;
411         dhp = dhp->dh_next;
412     }

```

```

414     /*
415      * Inform the vnode of the new mapping.
416      */
417     /*
418      * It is ok to use pass sdp->maxprot to ADDMAP rather than to use
419      * dhp specific maxprot because spec_addmap does not use maxprot.
420      */
421     error = VOP_ADDMAP(VTOCVP(sdp->vp), sdp->offset,
422         seg->s_as, seg->s_base, seg->s_size,
423         sdp->prot, sdp->maxprot, sdp->type, CRED(), NULL);

425     if (error != 0) {
426         sdp->devmap_data = NULL;
427         hat_unload(seg->s_as->a_hat, seg->s_base, seg->s_size,
428             HAT_UNLOAD_UNMAP);
429     } else {
430         /*
431          * Mappings of /dev/null don't count towards the VSZ of a
432          * process. Mappings of /dev/null have no mapping type.
433          */
434         if ((segop_gettype(seg, seg->s_base) & (MAP_SHARED |
435             if ((SEGOP_GETTYPE(seg, (seg)->s_base) & (MAP_SHARED |
436                 MAP_PRIVATE)) == 0) {
437                 seg->s_as->a_resvsize -= seg->s_size;
438             }
439         }

440     return (error);
441 }
_____unchanged_portion_omitted_____

```

new/usr/src/uts/common/vm/seg_dev.h

1

```
*****
4471 Tue Nov 24 09:34:48 2015
new/usr/src/uts/common/vm/seg_dev.h
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

113 #ifdef _KERNEL

115 /*
116  * Mappings of /dev/null come from segdev and have no mapping type.
117  */

119 #define SEG_IS_DEVNULL_MAPPING(seg) \
120     ((seg)->s_ops == &segdev_ops && \
121      ((segop_gettype((seg), (seg)->s_base) & \
122       (MAP_SHARED | MAP_PRIVATE)) == 0))
121     ((SEGOP_GETTYPE(seg, (seg)->s_base) & (MAP_SHARED | MAP_PRIVATE)) == 0))

124 extern void segdev_init(void);

126 extern int segdev_create(struct seg *, void *);

128 extern int segdev_copyto(struct seg *, caddr_t, const void *, void *, size_t);
129 extern int segdev_copyfrom(struct seg *, caddr_t, const void *, void *, size_t);
130 extern struct seg_ops segdev_ops;

132 #endif /* _KERNEL */

134 #ifdef __cplusplus
135 }
_____unchanged_portion_omitted_____
```

```

*****
45463 Tue Nov 24 09:34:48 2015
new/usr/src/uts/common/vm/seg_kmem.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

437 #define SEGMEM_BADOP(t)      (t(*)())segkmem_badop

439 /*ARGSUSED*/
440 static faultcode_t
441 segkmem_fault(struct hat *hat, struct seg *seg, caddr_t addr, size_t size,
442              enum fault_type type, enum seg_rw rw)
443 {
444     pgcnt_t npages;
445     spgcnt_t pg;
446     page_t *pp;
447     struct vnode *vp = seg->s_data;

449     ASSERT(RW_READ_HELD(&seg->s_as->a_lock));

451     if (seg->s_as != &kas || size > seg->s_size ||
452         addr < seg->s_base || addr + size > seg->s_base + seg->s_size)
453         panic("segkmem_fault: bad args");

455     /*
456      * If it is one of segkp pages, call segkp_fault.
457      */
458     if (segkp_bitmap && seg == &kvseg &&
459         BT_TEST(segkp_bitmap, btop((uintptr_t)(addr - seg->s_base))))
460         return (segop_fault(hat, segkp, addr, size, type, rw));
461         return (SEGOP_FAULT(hat, segkp, addr, size, type, rw));

462     if (rw != S_READ && rw != S_WRITE && rw != S_OTHER)
463         return (FC_NOSUPPORT);

465     npages = btopr(size);

467     switch (type) {
468     case F_SOFTLOCK:
469         /* lock down already-loaded translations */
470         for (pg = 0; pg < npages; pg++) {
471             pp = page_lookup(vp, (u_offset_t)(uintptr_t)addr,
472                             SE_SHARED);
473             if (pp == NULL) {
474                 /*
475                  * Hmm, no page. Does a kernel mapping
476                  * exist for it?
477                  */
478                 if (!hat_probe(kas.a_hat, addr)) {
479                     addr -= PAGE_SIZE;
480                     while (--pg >= 0) {
481                         pp = page_find(vp, (u_offset_t)
482                                       (uintptr_t)addr);
483                         if (pp)
484                             page_unlock(pp);
485                         addr -= PAGE_SIZE;
486                     }
487                     return (FC_NOMAP);
488                 }
489                 addr += PAGE_SIZE;
490             }
491             if (rw == S_OTHER)
492                 hat_reserve(seg->s_as, addr, size);
493             return (0);
494         case F_SOFTUNLOCK:

```

```

495         while (npages-- > 0) {
496             pp = page_find(vp, (u_offset_t)(uintptr_t)addr);
497             if (pp)
498                 page_unlock(pp);
499             addr += PAGE_SIZE;
500         }
501         return (0);
502     default:
503         return (FC_NOSUPPORT);
504     }
505     /*NOTREACHED*/
506 }

508 static int
509 segkmem_setprot(struct seg *seg, caddr_t addr, size_t size, uint_t prot)
510 {
511     ASSERT(RW_LOCK_HELD(&seg->s_as->a_lock));

513     if (seg->s_as != &kas || size > seg->s_size ||
514         addr < seg->s_base || addr + size > seg->s_base + seg->s_size)
515         panic("segkmem_setprot: bad args");

517     /*
518      * If it is one of segkp pages, call segkp.
519      */
520     if (segkp_bitmap && seg == &kvseg &&
521         BT_TEST(segkp_bitmap, btop((uintptr_t)(addr - seg->s_base))))
522         return (segop_setprot(segkp, addr, size, prot));
523         return (SEGOP_SETPROT(segkp, addr, size, prot));

524     if (prot == 0)
525         hat_unload(kas.a_hat, addr, size, HAT_UNLOAD);
526     else
527         hat_chgprot(kas.a_hat, addr, size, prot);
528     return (0);
529 }

531 /*
532 * This is a dummy segkmem function overloaded to call segkp
533 * when segkp is under the heap.
534 */
535 /* ARGSUSED */
536 static int
537 segkmem_checkprot(struct seg *seg, caddr_t addr, size_t size, uint_t prot)
538 {
539     ASSERT(RW_LOCK_HELD(&seg->s_as->a_lock));

541     if (seg->s_as != &kas)
542         segkmem_badop();

544     /*
545      * If it is one of segkp pages, call into segkp.
546      */
547     if (segkp_bitmap && seg == &kvseg &&
548         BT_TEST(segkp_bitmap, btop((uintptr_t)(addr - seg->s_base))))
549         return (segop_checkprot(segkp, addr, size, prot));
550         return (SEGOP_CHECKPROT(segkp, addr, size, prot));

551     segkmem_badop();
552     return (0);
553 }

555 /*
556 * This is a dummy segkmem function overloaded to call segkp
557 * when segkp is under the heap.
558 */

```

```

559 /* ARGSUSED */
560 static int
561 segkmem_kluster(struct seg *seg, caddr_t addr, ssize_t delta)
562 {
563     ASSERT(RW_LOCK_HELD(&seg->s_as->a_lock));
564
565     if (seg->s_as != &kas)
566         segkmem_badop();
567
568     /*
569     * If it is one of segkp pages, call into segkp.
570     */
571     if (segkp_bitmap && seg == &kvseg &&
572         BT_TEST(segkp_bitmap, btop((uintptr_t)(addr - seg->s_base))))
573         return (segop_kluster(segkp, addr, delta));
574     return (SEGOP_KLUSTER(segkp, addr, delta));
575 }
576
577 segkmem_badop();
578 return (0);
579 }
580
581 unchanged_portion_omitted
582
583 /*
584 * lock/unlock kmem pages over a given range [addr, addr+len).
585 * Returns a shadow list of pages in ppp. If there are holes
586 * in the range (e.g. some of the kernel mappings do not have
587 * underlying page_ts) returns ENOTSUP so that as_pagelock()
588 * will handle the range via as_fault(F_SOFTLOCK).
589 */
590 /*ARGSUSED*/
591 static int
592 segkmem_pagelock(struct seg *seg, caddr_t addr, size_t len,
593                 page_t ***ppp, enum lock_type type, enum seg_rw rw)
594 {
595     page_t **pplist, *pp;
596     pgcnt_t npages;
597     spgcnt_t pg;
598     size_t nb;
599     struct vnode *vp = seg->s_data;
600
601     ASSERT(ppp != NULL);
602
603     /*
604     * If it is one of segkp pages, call into segkp.
605     */
606     if (segkp_bitmap && seg == &kvseg &&
607         BT_TEST(segkp_bitmap, btop((uintptr_t)(addr - seg->s_base))))
608         return (segop_pagelock(segkp, addr, len, ppp, type, rw));
609     return (SEGOP_PAGELOCK(segkp, addr, len, ppp, type, rw));
610
611     npages = btopr(len);
612     nb = sizeof (page_t *) * npages;
613
614     if (type == L_PAGEUNLOCK) {
615         pplist = *ppp;
616         ASSERT(pplist != NULL);
617
618         for (pg = 0; pg < npages; pg++) {
619             pp = pplist[pg];
620             page_unlock(pp);
621         }
622         kmem_free(pplist, nb);
623         return (0);
624     }
625
626     ASSERT(type == L_PAGELOCK);

```

```

713     pplist = kmem_alloc(nb, KM_NOSLEEP);
714     if (pplist == NULL) {
715         *ppp = NULL;
716         return (ENOTSUP); /* take the slow path */
717     }
718
719     for (pg = 0; pg < npages; pg++) {
720         pp = page_lookup(vp, (u_offset_t)(uintptr_t)addr, SE_SHARED);
721         if (pp == NULL) {
722             while (--pg >= 0)
723                 page_unlock(pplist[pg]);
724             kmem_free(pplist, nb);
725             *ppp = NULL;
726             return (ENOTSUP);
727         }
728         pplist[pg] = pp;
729         addr += PAGE_SIZE;
730     }
731
732     *ppp = pplist;
733     return (0);
734 }
735
736 /*
737 * This is a dummy segkmem function overloaded to call segkp
738 * when segkp is under the heap.
739 */
740 /* ARGSUSED */
741 static int
742 segkmem_getmemid(struct seg *seg, caddr_t addr, memid_t *memidp)
743 {
744     ASSERT(RW_LOCK_HELD(&seg->s_as->a_lock));
745
746     if (seg->s_as != &kas)
747         segkmem_badop();
748
749     /*
750     * If it is one of segkp pages, call into segkp.
751     */
752     if (segkp_bitmap && seg == &kvseg &&
753         BT_TEST(segkp_bitmap, btop((uintptr_t)(addr - seg->s_base))))
754         return (segop_getmemid(segkp, addr, memidp));
755     return (SEGOP_GETMEMID(segkp, addr, memidp));
756
757     segkmem_badop();
758     return (0);
759 }
760
761 unchanged_portion_omitted

```

```

*****
286002 Tue Nov 24 09:34:48 2015
new/usr/src/uts/common/vm/seg_vn.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

6086 /*
6087  * segvn_setpagesize is called via segop_setpagesize from as_setpagesize,
6088  * segvn_setpagesize is called via SEGOP_SETPAGESIZE from as_setpagesize,
6089  * to determine if the seg is capable of mapping the requested szc.
6090  */
6091 static int
6092 segvn_setpagesize(struct seg *seg, caddr_t addr, size_t len, uint_t szc)
6093 {
6094     struct segvn_data *svd = (struct segvn_data *)seg->s_data;
6095     struct segvn_data *nsvd;
6096     struct anon_map *amp = svd->amp;
6097     struct seg *nseg;
6098     caddr_t eaddr = addr + len, a;
6099     size_t pgsz = page_get_pagesize(szc);
6100     pgcnt_t pgcnt = page_get_pagecnt(szc);
6101     int err;
6102     u_offset_t off = svd->offset + (uintptr_t)(addr - seg->s_base);

6103     ASSERT(seg->s_as && AS_WRITE_HELD(seg->s_as, &seg->s_as->a_lock));
6104     ASSERT(addr >= seg->s_base && eaddr <= seg->s_base + seg->s_size);

6106     if (seg->s_szc == szc || segvn_lpg_disable != 0) {
6107         return (0);
6108     }

6110     /*
6111      * addr should always be pgsz aligned but eaddr may be misaligned if
6112      * it's at the end of the segment.
6113      *
6114      * XXX we should assert this condition since as_setpagesize() logic
6115      * guarantees it.
6116      */
6117     if (!IS_P2ALIGNED(addr, pgsz) ||
6118         (!IS_P2ALIGNED(eaddr, pgsz) &&
6119          eaddr != seg->s_base + seg->s_size)) {

6121         segvn_setpgsz_align_err++;
6122         return (EINVAL);
6123     }

6125     if (amp != NULL && svd->type == MAP_SHARED) {
6126         ulong_t an_idx = svd->anon_index + seg_page(seg, addr);
6127         if (!IS_P2ALIGNED(an_idx, pgcnt)) {

6129             segvn_setpgsz_anon_align_err++;
6130             return (EINVAL);
6131         }
6132     }

6134     if ((svd->flags & MAP_NORESERVE) || seg->s_as == &kas ||
6135         szc > segvn_maxpgsz) {
6136         return (EINVAL);
6137     }

6139     /* paranoid check */
6140     if (svd->vp != NULL &&
6141         (IS_SWAPFSVP(svd->vp) || VN_ISKAS(svd->vp))) {
6142         return (EINVAL);
6143     }

```

```

6145     if (seg->s_szc == 0 && svd->vp != NULL &&
6146         map_addr_vacalign_check(addr, off)) {
6147         return (EINVAL);
6148     }

6150     /*
6151      * Check that protections are the same within new page
6152      * size boundaries.
6153      */
6154     if (svd->pageprot) {
6155         for (a = addr; a < eaddr; a += pgsz) {
6156             if ((a + pgsz) > eaddr) {
6157                 if (!sameprot(seg, a, eaddr - a)) {
6158                     return (EINVAL);
6159                 }
6160             } else {
6161                 if (!sameprot(seg, a, pgsz)) {
6162                     return (EINVAL);
6163                 }
6164             }
6165         }
6166     }

6168     /*
6169      * Since we are changing page size we first have to flush
6170      * the cache. This makes sure all the pagelock calls have
6171      * to recheck protections.
6172      */
6173     if (svd->softlockcnt > 0) {
6174         ASSERT(svd->tr_state == SEGVN_TR_OFF);

6176         /*
6177          * If this is shared segment non 0 softlockcnt
6178          * means locked pages are still in use.
6179          */
6180         if (svd->type == MAP_SHARED) {
6181             return (EAGAIN);
6182         }

6184         /*
6185          * Since we do have the segvn writers lock nobody can fill
6186          * the cache with entries belonging to this seg during
6187          * the purge. The flush either succeeds or we still have
6188          * pending I/Os.
6189          */
6190         segvn_purge(seg);
6191         if (svd->softlockcnt > 0) {
6192             return (EAGAIN);
6193         }
6194     }

6196     if (HAT_IS_REGION_COOKIE_VALID(svd->rcookie)) {
6197         ASSERT(svd->amp == NULL);
6198         ASSERT(svd->tr_state == SEGVN_TR_OFF);
6199         hat_leave_region(seg->s_as->a_hat, svd->rcookie,
6200             HAT_REGION_TEXT);
6201         svd->rcookie = HAT_INVALID_REGION_COOKIE;
6202     } else if (svd->tr_state == SEGVN_TR_INIT) {
6203         svd->tr_state = SEGVN_TR_OFF;
6204     } else if (svd->tr_state == SEGVN_TR_ON) {
6205         ASSERT(svd->amp != NULL);
6206         segvn_textunrepl(seg, 1);
6207         ASSERT(svd->amp == NULL && svd->tr_state == SEGVN_TR_OFF);
6208         amp = NULL;
6209     }

```



```

6211     /*
6212     * Operation for sub range of existing segment.
6213     */
6214     if (addr != seg->s_base || eaddr != (seg->s_base + seg->s_size)) {
6215         if (szc < seg->s_szc) {
6216             VM_STAT_ADD(segvmstats.demoterange[2]);
6217             err = segvn_demote_range(seg, addr, len, SDR_RANGE, 0);
6218             if (err == 0) {
6219                 return (IE_RETRY);
6220             }
6221             if (err == ENOMEM) {
6222                 return (IE_NOMEM);
6223             }
6224             return (err);
6225         }
6226         if (addr != seg->s_base) {
6227             nseg = segvn_split_seg(seg, addr);
6228             if (eaddr != (nseg->s_base + nseg->s_size)) {
6229                 /* eaddr is szc aligned */
6230                 (void) segvn_split_seg(nseg, eaddr);
6231             }
6232             return (IE_RETRY);
6233         }
6234         if (eaddr != (seg->s_base + seg->s_size)) {
6235             /* eaddr is szc aligned */
6236             (void) segvn_split_seg(seg, eaddr);
6237         }
6238         return (IE_RETRY);
6239     }

6241     /*
6242     * Break any low level sharing and reset seg->s_szc to 0.
6243     */
6244     if ((err = segvn_clrsrc(seg)) != 0) {
6245         if (err == ENOMEM) {
6246             err = IE_NOMEM;
6247         }
6248         return (err);
6249     }
6250     ASSERT(seg->s_szc == 0);

6252     /*
6253     * If the end of the current segment is not pgsz aligned
6254     * then attempt to concatenate with the next segment.
6255     */
6256     if (!IS_P2ALIGNED(eaddr, pgsz)) {
6257         nseg = AS_SEGNEXT(seg->s_as, seg);
6258         if (nseg == NULL || nseg == seg || eaddr != nseg->s_base) {
6259             return (ENOMEM);
6260         }
6261         if (nseg->s_ops != &segvn_ops) {
6262             return (EINVAL);
6263         }
6264         nsvd = (struct segvn_data *)nseg->s_data;
6265         if (nsvd->softlockcnt > 0) {
6266             /*
6267              * If this is shared segment non 0 softlockcnt
6268              * means locked pages are still in use.
6269              */
6270             if (nsvd->type == MAP_SHARED) {
6271                 return (EAGAIN);
6272             }
6273             segvn_purge(nseg);
6274             if (nsvd->softlockcnt > 0) {
6275                 return (EAGAIN);

```

```

6276         }
6277     }
6278     err = segvn_clrsrc(nseg);
6279     if (err == ENOMEM) {
6280         err = IE_NOMEM;
6281     }
6282     if (err != 0) {
6283         return (err);
6284     }
6285     ASSERT(nsvd->rcookie == HAT_INVALID_REGION_COOKIE);
6286     err = segvn_concat(seg, nseg, 1);
6287     if (err == -1) {
6288         return (EINVAL);
6289     }
6290     if (err == -2) {
6291         return (IE_NOMEM);
6292     }
6293     return (IE_RETRY);
6294 }

6296     /*
6297     * May need to re-align anon array to
6298     * new szc.
6299     */
6300     if (amp != NULL) {
6301         if (!IS_P2ALIGNED(svd->anon_index, pgcnt)) {
6302             struct anon_hdr *nahp;

6304             ASSERT(svd->type == MAP_PRIVATE);

6306             ANON_LOCK_ENTER(&amp->a_rwlock, RW_WRITER);
6307             ASSERT(amp->refcnt == 1);
6308             nahp = anon_create(btop(amp->size), ANON_NOSLEEP);
6309             if (nahp == NULL) {
6310                 ANON_LOCK_EXIT(&amp->a_rwlock);
6311                 return (IE_NOMEM);
6312             }
6313             if (anon_copy_ptr(amp->ahp, svd->anon_index,
6314                 nahp, 0, btop(seg->s_size), ANON_NOSLEEP)) {
6315                 anon_release(nahp, btop(amp->size));
6316                 ANON_LOCK_EXIT(&amp->a_rwlock);
6317                 return (IE_NOMEM);
6318             }
6319             anon_release(amp->ahp, btop(amp->size));
6320             amp->ahp = nahp;
6321             svd->anon_index = 0;
6322             ANON_LOCK_EXIT(&amp->a_rwlock);
6323         }
6324     }
6325     if (svd->vp != NULL && szc != 0) {
6326         struct vattr va;
6327         u_offset_t eoffpage = svd->offset;
6328         va.va_mask = AT_SIZE;
6329         eoffpage += seg->s_size;
6330         eoffpage = btop(eoffpage);
6331         if (VOP_GETATTR(svd->vp, &va, 0, svd->cred, NULL) != 0) {
6332             segvn_setpgsz_getattr_err++;
6333             return (EINVAL);
6334         }
6335         if (btop(va.va_size) < eoffpage) {
6336             segvn_setpgsz_eof_err++;
6337             return (EINVAL);
6338         }
6339         if (amp != NULL) {
6340             /*
6341              * anon_fill_cow_holes() may call VOP_GETPAGE().

```

```
6342         * don't take anon map lock here to avoid holding it
6343         * across VOP_GETPAGE() calls that may call back into
6344         * segvn for klsutering checks. We don't really need
6345         * anon map lock here since it's a private segment and
6346         * we hold as level lock as writers.
6347         */
6348         if ((err = anon_fill_cow_holes(seg, seg->s_base,
6349             amp->ahp, svd->anon_index, svd->vp, svd->offset,
6350             seg->s_size, szc, svd->prot, svd->vpage,
6351             svd->cred)) != 0) {
6352             return (EINVAL);
6353         }
6354     }
6355     segvn_setvnode_mpss(svd->vp);
6356 }
6357
6358 if (amp != NULL) {
6359     ANON_LOCK_ENTER(&amp->a_rwlock, RW_WRITER);
6360     if (svd->type == MAP_PRIVATE) {
6361         amp->a_szc = szc;
6362     } else if (szc > amp->a_szc) {
6363         amp->a_szc = szc;
6364     }
6365     ANON_LOCK_EXIT(&amp->a_rwlock);
6366 }
6367
6368     seg->s_szc = szc;
6369
6370     return (0);
6371 }
```

unchanged portion omitted

```

*****
94702 Tue Nov 24 09:34:49 2015
new/usr/src/uts/common/vm/vm_as.c
patch lower-case-segops
*****
_____unchanged_portion_omitted_____

679 /*
680  * Free an address space data structure.
681  * Need to free the hat first and then
682  * all the segments on this as and finally
683  * the space for the as struct itself.
684  */
685 void
686 as_free(struct as *as)
687 {
688     struct hat *hat = as->a_hat;
689     struct seg *seg, *next;
690     int called = 0;

692 top:
693     /*
694     * Invoke ALL callbacks. as_do_callbacks will do one callback
695     * per call, and not return (-1) until the callback has completed.
696     * When as_do_callbacks returns zero, all callbacks have completed.
697     */
698     mutex_enter(&as->a_contents);
699     while (as->a_callbacks && as_do_callbacks(as, AS_ALL_EVENT, 0, 0))
700         ;

702     /* This will prevent new XHATs from attaching to as */
703     if (!called)
704         AS_SETBUSY(as);
705     mutex_exit(&as->a_contents);
706     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);

708     if (!called) {
709         called = 1;
710         hat_free_start(hat);
711         if (as->a_xhat != NULL)
712             xhat_free_start_all(as);
713     }
714     for (seg = AS_SEGFIRST(as); seg != NULL; seg = next) {
715         int err;

717         next = AS_SEGNEXT(as, seg);
718     retry:
719         err = segop_unmap(seg, seg->s_base, seg->s_size);
720         err = SEGOP_UNMAP(seg, seg->s_base, seg->s_size);
721         if (err == EAGAIN) {
722             mutex_enter(&as->a_contents);
723             if (as->a_callbacks) {
724                 AS_LOCK_EXIT(as, &as->a_lock);
725             } else if (!AS_ISNOUNMAPWAIT(as)) {
726                 /*
727                  * Memory is currently locked. Wait for a
728                  * cv_signal that it has been unlocked, then
729                  * try the operation again.
730                  */
731                 if (AS_ISUNMAPWAIT(as) == 0)
732                     cv_broadcast(&as->a_cv);
733                 AS_SETUNMAPWAIT(as);
734                 AS_LOCK_EXIT(as, &as->a_lock);
735                 while (AS_ISUNMAPWAIT(as))
736                     cv_wait(&as->a_cv, &as->a_contents);
737             } else {

```

```

737     /*
738     * We may have raced with
739     * segvn_reclaim()/segspt_reclaim(). In this
740     * case clean nounmapwait flag and retry since
741     * softlockcnt in this segment may be already
742     * 0. We don't drop as writer lock so our
743     * number of retries without sleeping should
744     * be very small. See segvn_reclaim() for
745     * more comments.
746     */
747     AS_CLRNOUNMAPWAIT(as);
748     mutex_exit(&as->a_contents);
749     goto retry;
750 }
751     mutex_exit(&as->a_contents);
752     goto top;
753 } else {
754     /*
755     * We do not expect any other error return at this
756     * time. This is similar to an ASSERT in seg_unmap()
757     */
758     ASSERT(err == 0);
759 }
760 }
761 hat_free_end(hat);
762 if (as->a_xhat != NULL)
763     xhat_free_end_all(as);
764 AS_LOCK_EXIT(as, &as->a_lock);

766 /* /proc stuff */
767 ASSERT(avl_numnodes(&as->a_wpage) == 0);
768 if (as->a_objectdir) {
769     kmem_free(as->a_objectdir, as->a_sizedir * sizeof (vnode_t *));
770     as->a_objectdir = NULL;
771     as->a_sizedir = 0;
772 }

774 /*
775  * Free the struct as back to kmem. Assert it has no segments.
776  */
777 ASSERT(avl_numnodes(&as->a_segtree) == 0);
778 kmem_cache_free(as_cache, as);
779 }

781 int
782 as_dup(struct as *as, struct proc *forkedproc)
783 {
784     struct as *newas;
785     struct seg *seg, *newseg;
786     size_t purgesize = 0;
787     int error;

789     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
790     as_clearwatch(as);
791     newas = as_alloc();
792     newas->a_userlimit = as->a_userlimit;
793     newas->a_proc = forkedproc;

795     AS_LOCK_ENTER(newas, &newas->a_lock, RW_WRITER);

797     /* This will prevent new XHATs from attaching */
798     mutex_enter(&as->a_contents);
799     AS_SETBUSY(as);
800     mutex_exit(&as->a_contents);
801     mutex_enter(&newas->a_contents);
802     AS_SETBUSY(newas);

```

```

803     mutex_exit(&newas->a_contents);
805     (void) hat_dup(as->a_hat, newas->a_hat, NULL, 0, HAT_DUP_SRD);
807     for (seg = AS_SEGFIRST(as); seg != NULL; seg = AS_SEGNEXT(as, seg)) {
809         if (seg->s_flags & S_PURGE) {
810             purgesize += seg->s_size;
811             continue;
812         }
814         newseg = seg_alloc(newas, seg->s_base, seg->s_size);
815         if (newseg == NULL) {
816             AS_LOCK_EXIT(newas, &newas->a_lock);
817             as_setwatch(as);
818             mutex_enter(&as->a_contents);
819             AS_CLRBUSY(as);
820             mutex_exit(&as->a_contents);
821             AS_LOCK_EXIT(as, &as->a_lock);
822             as_free(newas);
823             return (-1);
824         }
825         if ((error = segop_dup(seg, newseg)) != 0) {
826             if ((error = SEGOP_DUP(seg, newseg)) != 0) {
827                 /*
828                  * We call seg_free() on the new seg
829                  * because the segment is not set up
830                  * completely; i.e. it has no ops.
831                  */
832                 as_setwatch(as);
833                 mutex_enter(&as->a_contents);
834                 AS_CLRBUSY(as);
835                 mutex_exit(&as->a_contents);
836                 AS_LOCK_EXIT(as, &as->a_lock);
837                 seg_free(newseg);
838                 AS_LOCK_EXIT(newas, &newas->a_lock);
839                 as_free(newas);
840                 return (error);
841             }
842             newas->a_size += seg->s_size;
843         }
844         newas->a_resvsize = as->a_resvsize - purgesize;
845     }
846     error = hat_dup(as->a_hat, newas->a_hat, NULL, 0, HAT_DUP_ALL);
847     if (as->a_xhat != NULL)
848         error |= xhat_dup_all(as, newas, NULL, 0, HAT_DUP_ALL);
849     mutex_enter(&newas->a_contents);
850     AS_CLRBUSY(newas);
851     mutex_exit(&newas->a_contents);
852     AS_LOCK_EXIT(newas, &newas->a_lock);
854     as_setwatch(as);
855     mutex_enter(&as->a_contents);
856     AS_CLRBUSY(as);
857     mutex_exit(&as->a_contents);
858     AS_LOCK_EXIT(as, &as->a_lock);
859     if (error != 0) {
860         as_free(newas);
861         return (error);
862     }
863     forkedproc->p_as = newas;
864     return (0);
865 }
867 /*

```

```

868 * Handle a ``fault`` at addr for size bytes.
869 */
870 faultcode_t
871 as_fault(struct hat *hat, struct as *as, caddr_t addr, size_t size,
872          enum fault_type type, enum seg_rw rw)
873 {
874     struct seg *seg;
875     caddr_t raddr; /* rounded down addr */
876     size_t rsize; /* rounded up size */
877     size_t ssize;
878     faultcode_t res = 0;
879     caddr_t addrsav;
880     struct seg *segsav;
881     int as_lock_held;
882     klpw_t *lwp = ttolwp(curthread);
883     int is_xhat = 0;
884     int holding_wpage = 0;
885     extern struct seg_ops segdev_ops;
889     if (as->a_hat != hat) {
890         /* This must be an XHAT then */
891         is_xhat = 1;
893         if ((type != F_INVALID) || (as == &kas))
894             return (FC_NOSUPPORT);
895     }
897     retry:
898     if (!is_xhat) {
899         /*
900          * Indicate that the lwp is not to be stopped while waiting
901          * for a pagefault. This is to avoid deadlock while debugging
902          * a process via /proc over NFS (in particular).
903          */
904         if (lwp != NULL)
905             lwp->lwp_nostop++;
907         /*
908          * same length must be used when we softlock and softunlock.
909          * We don't support softunlocking lengths less than
910          * the original length when there is largepage support.
911          * See seg_dev.c for more comments.
912          */
913         switch (type) {
915             case F_SOFTLOCK:
916                 CPU_STATS_ADD_K(vm, softlock, 1);
917                 break;
919             case F_SOFTUNLOCK:
920                 break;
922             case F_PROT:
923                 CPU_STATS_ADD_K(vm, prot_fault, 1);
924                 break;
926             case F_INVALID:
927                 CPU_STATS_ENTER_K();
928                 CPU_STATS_ADDQ(CPU, vm, as_fault, 1);
929                 if (as == &kas)
930                     CPU_STATS_ADDQ(CPU, vm, kernel_asflt, 1);
931                 CPU_STATS_EXIT_K();
932                 break;
933         }

```

```

934     }
935
936     /* Kernel probe */
937     TNF_PROBE_3(address_fault, "vm pagefault", /* CSTYLE */ ,
938               tnfn_opaque, address,      addr,
939               tnfn_fault_type, fault_type, type,
940               tnfn_seg_access, access,    rw);
941
942     raddr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);
943     rsize = (((size_t)(addr + size) + PAGEOFFSET) & PAGEMASK) -
944             (size_t)raddr;
945
946     /*
947     * XXX -- Don't grab the as lock for segkmap. We should grab it for
948     * correctness, but then we could be stuck holding this lock for
949     * a LONG time if the fault needs to be resolved on a slow
950     * filesystem, and then no-one will be able to exec new commands,
951     * as exec'ing requires the write lock on the as.
952     */
953     if (as == &kas && segkmap && segkmap->s_base <= raddr &&
954         raddr + size < segkmap->s_base + segkmap->s_size) {
955         /*
956         * if (as==&kas), this can't be XHAT: we've already returned
957         * FC_NOSUPPORT.
958         */
959         seg = segkmap;
960         as_lock_held = 0;
961     } else {
962         AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
963         if (is_xhat && avl_numnodes(&as->a_wpage) != 0) {
964             /*
965             * Grab and hold the writers' lock on the as
966             * if the fault is to a watched page.
967             * This will keep CPUs from "peeking" at the
968             * address range while we're temporarily boosting
969             * the permissions for the XHAT device to
970             * resolve the fault in the segment layer.
971             *
972             * We could check whether faulted address
973             * is within a watched page and only then grab
974             * the writer lock, but this is simpler.
975             */
976             AS_LOCK_EXIT(as, &as->a_lock);
977             AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
978         }
979
980         seg = as_segat(as, raddr);
981         if (seg == NULL) {
982             AS_LOCK_EXIT(as, &as->a_lock);
983             if ((lwp != NULL) && (!is_xhat))
984                 lwp->lwp_nostop--;
985             return (FC_NOMAP);
986         }
987
988         as_lock_held = 1;
989     }
990
991     addrsav = raddr;
992     segsav = seg;
993
994     for (; rsize != 0; rsize -= ssize, raddr += ssize) {
995         if (raddr >= seg->s_base + seg->s_size) {
996             seg = AS_SEGNEXT(as, seg);
997             if (seg == NULL || raddr != seg->s_base) {
998                 res = FC_NOMAP;
999                 break;

```

```

1000     }
1001     }
1002     if (raddr + rsize > seg->s_base + seg->s_size)
1003         ssize = seg->s_base + seg->s_size - raddr;
1004     else
1005         ssize = rsize;
1006
1007     if (!is_xhat || (seg->s_ops != &segdev_ops)) {
1008
1009         if (is_xhat && avl_numnodes(&as->a_wpage) != 0 &&
1010             pr_is_watchpage_as(raddr, rw, as)) {
1011             /*
1012             * Handle watch pages. If we're faulting on a
1013             * watched page from an X-hat, we have to
1014             * restore the original permissions while we
1015             * handle the fault.
1016             */
1017             as_clearwatch(as);
1018             holding_wpage = 1;
1019         }
1020
1021         res = segop_fault(hat, seg, raddr, ssize, type, rw);
1022         res = SEGOP_FAULT(hat, seg, raddr, ssize, type, rw);
1023
1024         /* Restore watchpoints */
1025         if (holding_wpage) {
1026             as_setwatch(as);
1027             holding_wpage = 0;
1028         }
1029
1030         if (res != 0)
1031             break;
1032     } else {
1033         /* XHAT does not support seg_dev */
1034         res = FC_NOSUPPORT;
1035         break;
1036     }
1037 }
1038
1039 /*
1040 * If we were SOFTLOCKING and encountered a failure,
1041 * we must SOFTUNLOCK the range we already did. (Maybe we
1042 * should just panic if we are SOFTLOCKING or even SOFTUNLOCKING
1043 * right here...)
1044 */
1045 if (res != 0 && type == F_SOFTLOCK) {
1046     for (seg = segsav; addrsav < raddr; addrsav += ssize) {
1047         if (addrsav >= seg->s_base + seg->s_size)
1048             seg = AS_SEGNEXT(as, seg);
1049         ASSERT(seg != NULL);
1050         /*
1051         * Now call the fault routine again to perform the
1052         * unlock using S_OTHER instead of the rw variable
1053         * since we never got a chance to touch the pages.
1054         */
1055         if (raddr > seg->s_base + seg->s_size)
1056             ssize = seg->s_base + seg->s_size - addrsav;
1057         else
1058             ssize = raddr - addrsav;
1059         (void) segop_fault(hat, seg, addrsav, ssize,
1060                          (void) SEGOP_FAULT(hat, seg, addrsav, ssize,
1061                                               F_SOFTUNLOCK, S_OTHER));
1062     }
1063     if (as_lock_held)
1064         AS_LOCK_EXIT(as, &as->a_lock);

```

```

1064     if ((lwp != NULL) && (!is_xhat))
1065         lwp->lwp_nostop--;

1067 /*
1068  * If the lower levels returned EDEADLK for a fault,
1069  * It means that we should retry the fault. Let's wait
1070  * a bit also to let the deadlock causing condition clear.
1071  * This is part of a gross hack to work around a design flaw
1072  * in the ufs/sds logging code and should go away when the
1073  * logging code is re-designed to fix the problem. See bug
1074  * 4125102 for details of the problem.
1075  */
1076     if (FC_ERRNO(res) == EDEADLK) {
1077         delay(deadlk_wait);
1078         res = 0;
1079         goto retry;
1080     }
1081     return (res);
1082 }

1086 /*
1087  * Asynchronous ``fault'' at addr for size bytes.
1088  */
1089     faultcode_t
1090     as_faulta(struct as *as, caddr_t addr, size_t size)
1091     {
1092         struct seg *seg;
1093         caddr_t raddr;                /* rounded down addr */
1094         size_t rsize;                /* rounded up size */
1095         faultcode_t res = 0;
1096         klwp_t *lwp = ttolwp(curthread);

1098     retry:
1099         /*
1100          * Indicate that the lwp is not to be stopped while waiting
1101          * for a pagefault. This is to avoid deadlock while debugging
1102          * a process via /proc over NFS (in particular).
1103          */
1104         if (lwp != NULL)
1105             lwp->lwp_nostop++;

1107         raddr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);
1108         rsize = (((size_t)(addr + size) + PAGEOFFSET) & PAGEMASK) -
1109             (size_t)raddr;

1111         AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
1112         seg = as_segat(as, raddr);
1113         if (seg == NULL) {
1114             AS_LOCK_EXIT(as, &as->a_lock);
1115             if (lwp != NULL)
1116                 lwp->lwp_nostop--;
1117             return (FC_NOMAP);
1118         }

1120         for (; rsize != 0; rsize -= PAGE_SIZE, raddr += PAGE_SIZE) {
1121             if (raddr >= seg->s_base + seg->s_size) {
1122                 seg = AS_SEGNEXT(as, seg);
1123                 if (seg == NULL || raddr != seg->s_base) {
1124                     res = FC_NOMAP;
1125                     break;
1126                 }
1127             }
1128             res = segop_faulta(seg, raddr);
1129             res = SEGOP_FAULTA(seg, raddr);

```

```

1129         if (res != 0)
1130             break;
1131     }
1132     AS_LOCK_EXIT(as, &as->a_lock);
1133     if (lwp != NULL)
1134         lwp->lwp_nostop--;
1135     /*
1136      * If the lower levels returned EDEADLK for a fault,
1137      * It means that we should retry the fault. Let's wait
1138      * a bit also to let the deadlock causing condition clear.
1139      * This is part of a gross hack to work around a design flaw
1140      * in the ufs/sds logging code and should go away when the
1141      * logging code is re-designed to fix the problem. See bug
1142      * 4125102 for details of the problem.
1143      */
1144     if (FC_ERRNO(res) == EDEADLK) {
1145         delay(deadlk_wait);
1146         res = 0;
1147         goto retry;
1148     }
1149     return (res);
1150 }

1152 /*
1153  * Set the virtual mapping for the interval from [addr : addr + size)
1154  * in address space 'as' to have the specified protection.
1155  * It is ok for the range to cross over several segments,
1156  * as long as they are contiguous.
1157  */
1158     int
1159     as_setprot(struct as *as, caddr_t addr, size_t size, uint_t prot)
1160     {
1161         struct seg *seg;
1162         struct as_callback *cb;
1163         size_t ssize;
1164         caddr_t raddr;                /* rounded down addr */
1165         size_t rsize;                /* rounded up size */
1166         int error = 0, writer = 0;
1167         caddr_t saveraddr;
1168         size_t saversize;

1170     setprot_top:
1171         raddr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);
1172         rsize = (((size_t)(addr + size) + PAGEOFFSET) & PAGEMASK) -
1173             (size_t)raddr;

1175         if (raddr + rsize < raddr)                /* check for wraparound */
1176             return (ENOMEM);

1178         saveraddr = raddr;
1179         saversize = rsize;

1181         /*
1182          * Normally we only lock the as as a reader. But
1183          * if due to setprot the segment driver needs to split
1184          * a segment it will return IE_RETRY. Therefore we re-acquire
1185          * the as lock as a writer so the segment driver can change
1186          * the seg list. Also the segment driver will return IE_RETRY
1187          * after it has changed the segment list so we therefore keep
1188          * locking as a writer. Since these operations should be rare
1189          * want to only lock as a writer when necessary.
1190          */
1191         if (writer || avl_numnodes(&as->a_wpage) != 0) {
1192             AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
1193         } else {
1194             AS_LOCK_ENTER(as, &as->a_lock, RW_READER);

```

```

1195     }
1197     as_clearwatchprot(as, raddr, rsize);
1198     seg = as_segat(as, raddr);
1199     if (seg == NULL) {
1200         as_setwatch(as);
1201         AS_LOCK_EXIT(as, &as->a_lock);
1202         return (ENOMEM);
1203     }
1205     for (; rsize != 0; rsize -= ssize, raddr += ssize) {
1206         if (raddr >= seg->s_base + seg->s_size) {
1207             seg = AS_SEGNEXT(as, seg);
1208             if (seg == NULL || raddr != seg->s_base) {
1209                 error = ENOMEM;
1210                 break;
1211             }
1212         }
1213         if ((raddr + rsize) > (seg->s_base + seg->s_size))
1214             ssize = seg->s_base + seg->s_size - raddr;
1215         else
1216             ssize = rsize;
1217     retry:
1218         error = segop_setprot(seg, raddr, ssize, prot);
1219         error = SEGOP_SETPROT(seg, raddr, ssize, prot);
1220
1221         if (error == IE_NOMEM) {
1222             error = EAGAIN;
1223             break;
1224         }
1225         if (error == IE_RETRY) {
1226             AS_LOCK_EXIT(as, &as->a_lock);
1227             writer = 1;
1228             goto setprot_top;
1229         }
1231         if (error == EAGAIN) {
1232             /*
1233              * Make sure we have a_lock as writer.
1234              */
1235             if (writer == 0) {
1236                 AS_LOCK_EXIT(as, &as->a_lock);
1237                 writer = 1;
1238                 goto setprot_top;
1239             }
1241             /*
1242              * Memory is currently locked. It must be unlocked
1243              * before this operation can succeed through a retry.
1244              * The possible reasons for locked memory and
1245              * corresponding strategies for unlocking are:
1246              * (1) Normal I/O
1247              *     wait for a signal that the I/O operation
1248              *     has completed and the memory is unlocked.
1249              * (2) Asynchronous I/O
1250              *     The aio subsystem does not unlock pages when
1251              *     the I/O is completed. Those pages are unlocked
1252              *     when the application calls aiowait/aioerror.
1253              *     So, to prevent blocking forever, cv_broadcast()
1254              *     is done to wake up aio_cleanup_thread.
1255              *     Subsequently, segvn_reclaim will be called, and
1256              *     that will do AS_CLRUNMAPWAIT() and wake us up.
1257              * (3) Long term page locking:
1258              *     Drivers intending to have pages locked for a
1259              *     period considerably longer than for normal I/O

```

```

1260         * (essentially forever) may have registered for a
1261         * callback so they may unlock these pages on
1262         * request. This is needed to allow this operation
1263         * to succeed. Each entry on the callback list is
1264         * examined. If the event or address range pertains
1265         * the callback is invoked (unless it already is in
1266         * progress). The a_contents lock must be dropped
1267         * before the callback, so only one callback can
1268         * be done at a time. Go to the top and do more
1269         * until zero is returned. If zero is returned,
1270         * either there were no callbacks for this event
1271         * or they were already in progress.
1272         */
1273         mutex_enter(&as->a_contents);
1274         if (as->a_callbacks &&
1275             (cb = as_find_callback(as, AS_SETPROT_EVENT,
1276                 seg->s_base, seg->s_size))) {
1277             AS_LOCK_EXIT(as, &as->a_lock);
1278             as_execute_callback(as, cb, AS_SETPROT_EVENT);
1279         } else if (!AS_ISNOUNMAPWAIT(as)) {
1280             if (AS_ISUNMAPWAIT(as) == 0)
1281                 cv_broadcast(&as->a_cv);
1282             AS_SETUNMAPWAIT(as);
1283             AS_LOCK_EXIT(as, &as->a_lock);
1284             while (AS_ISUNMAPWAIT(as))
1285                 cv_wait(&as->a_cv, &as->a_contents);
1286         } else {
1287             /*
1288              * We may have raced with
1289              * segvn_reclaim()/segspt_reclaim(). In this
1290              * case clean nounmapwait flag and retry since
1291              * softlocknt in this segment may be already
1292              * 0. We don't drop as writer lock so our
1293              * number of retries without sleeping should
1294              * be very small. See segvn_reclaim() for
1295              * more comments.
1296              */
1297             AS_CLRNOUNMAPWAIT(as);
1298             mutex_exit(&as->a_contents);
1299             goto retry;
1300         }
1301         mutex_exit(&as->a_contents);
1302         goto setprot_top;
1303     } else if (error != 0)
1304         break;
1305     }
1306     if (error != 0) {
1307         as_setwatch(as);
1308     } else {
1309         as_setwatchprot(as, saveraddr, saversize, prot);
1310     }
1311     AS_LOCK_EXIT(as, &as->a_lock);
1312     return (error);
1313 }
1315 /*
1316  * Check to make sure that the interval [addr, addr + size)
1317  * in address space 'as' has at least the specified protection.
1318  * It is ok for the range to cross over several segments, as long
1319  * as they are contiguous.
1320  */
1321 int
1322 as_checkprot(struct as *as, caddr_t addr, size_t size, uint_t prot)
1323 {
1324     struct seg *seg;
1325     size_t ssize;

```

```

1326     caddr_t raddr;                /* rounded down addr */
1327     size_t rsize;                 /* rounded up size */
1328     int error = 0;

1330     raddr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);
1331     rsize = (((size_t)(addr + size) + PAGEOFFSET) & PAGEMASK) -
1332         (size_t)raddr;

1334     if (raddr + rsize < raddr)    /* check for wraparound */
1335         return (ENOMEM);

1337     /*
1338     * This is ugly as sin...
1339     * Normally, we only acquire the address space readers lock.
1340     * However, if the address space has watchpoints present,
1341     * we must acquire the writer lock on the address space for
1342     * the benefit of as_clearwatchprot() and as_setwatchprot().
1343     */
1344     if (avl_numnodes(&as->a_wpage) != 0)
1345         AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
1346     else
1347         AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
1348     as_clearwatchprot(as, raddr, rsize);
1349     seg = as_segat(as, raddr);
1350     if (seg == NULL) {
1351         as_setwatch(as);
1352         AS_LOCK_EXIT(as, &as->a_lock);
1353         return (ENOMEM);
1354     }

1356     for (; rsize != 0; rsize -= ssize, raddr += ssize) {
1357         if (raddr >= seg->s_base + seg->s_size) {
1358             seg = AS_SEGNEXT(as, seg);
1359             if (seg == NULL || raddr != seg->s_base) {
1360                 error = ENOMEM;
1361                 break;
1362             }
1363         }
1364         if ((raddr + rsize) > (seg->s_base + seg->s_size))
1365             ssize = seg->s_base + seg->s_size - raddr;
1366         else
1367             ssize = rsize;

1369         error = segop_checkprot(seg, raddr, ssize, prot);
1370         error = SEGOP_CHECKPROT(seg, raddr, ssize, prot);
1371         if (error != 0)
1372             break;
1373     }
1374     as_setwatch(as);
1375     AS_LOCK_EXIT(as, &as->a_lock);
1376     return (error);
1377 }

1378 int
1379 as_unmap(struct as *as, caddr_t addr, size_t size)
1380 {
1381     struct seg *seg, *seg_next;
1382     struct as_callback *cb;
1383     caddr_t raddr, eaddr;
1384     size_t ssize, rsize = 0;
1385     int err;

1387 top:
1388     raddr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);
1389     eaddr = (caddr_t)((uintptr_t)(addr + size) + PAGEOFFSET) &
1390         (uintptr_t)PAGEMASK);

```

```

1392     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);

1394     as->a_updatedir = 1;          /* inform /proc */
1395     gethrstime(&as->a_updatetime);

1397     /*
1398     * Use as_findseg to find the first segment in the range, then
1399     * step through the segments in order, following s_next.
1400     */
1401     as_clearwatchprot(as, raddr, eaddr - raddr);

1403     for (seg = as_findseg(as, raddr, 0); seg != NULL; seg = seg_next) {
1404         if (eaddr <= seg->s_base)
1405             break;                /* eaddr was in a gap; all done */

1407         /* this is implied by the test above */
1408         ASSERT(raddr < eaddr);

1410         if (raddr < seg->s_base)
1411             raddr = seg->s_base;    /* raddr was in a gap */

1413         if (eaddr > (seg->s_base + seg->s_size))
1414             ssize = seg->s_base + seg->s_size - raddr;
1415         else
1416             ssize = eaddr - raddr;

1418         /*
1419         * Save next segment pointer since seg can be
1420         * destroyed during the segment unmap operation.
1421         */
1422         seg_next = AS_SEGNEXT(as, seg);

1424         /*
1425         * We didn't count /dev/null mappings, so ignore them here.
1426         * We'll handle MAP_NORESERVE cases in segvn_unmap(). (Again,
1427         * we have to do this check here while we have seg.)
1428         */
1429         rsize = 0;
1430         if (!SEG_IS_DEVNULL_MAPPING(seg) &&
1431             !SEG_IS_PARTIAL_RESV(seg))
1432             rsize = ssize;

1434     retry:
1435     err = segop_unmap(seg, raddr, ssize);
1436     err = SEGOP_UNMAP(seg, raddr, ssize);
1437     if (err == EAGAIN) {
1438         /*
1439         * Memory is currently locked. It must be unlocked
1440         * before this operation can succeed through a retry.
1441         * The possible reasons for locked memory and
1442         * corresponding strategies for unlocking are:
1443         * (1) Normal I/O
1444         *     wait for a signal that the I/O operation
1445         *     has completed and the memory is unlocked.
1446         * (2) Asynchronous I/O
1447         *     The aio subsystem does not unlock pages when
1448         *     the I/O is completed. Those pages are unlocked
1449         *     when the application calls aiowait/aioerror.
1450         *     So, to prevent blocking forever, cv_broadcast()
1451         *     is done to wake up aio_cleanup_thread.
1452         *     Subsequently, segvn_reclaim will be called, and
1453         *     that will do AS_CLRUNMAPWAIT() and wake us up.
1454         * (3) Long term page locking:
1455         *     Drivers intending to have pages locked for a
1456         *     period considerably longer than for normal I/O

```



```

1456      * (essentially forever) may have registered for a
1457      * callback so they may unlock these pages on
1458      * request. This is needed to allow this operation
1459      * to succeed. Each entry on the callback list is
1460      * examined. If the event or address range pertains
1461      * the callback is invoked (unless it already is in
1462      * progress). The a_contents lock must be dropped
1463      * before the callback, so only one callback can
1464      * be done at a time. Go to the top and do more
1465      * until zero is returned. If zero is returned,
1466      * either there were no callbacks for this event
1467      * or they were already in progress.
1468      */
1469     mutex_enter(&as->a_contents);
1470     if (as->a_callbacks &&
1471         (cb = as_find_callback(as, AS_UNMAP_EVENT,
1472             seg->s_base, seg->s_size))) {
1473         AS_LOCK_EXIT(as, &as->a_lock);
1474         as_execute_callback(as, cb, AS_UNMAP_EVENT);
1475     } else if (!AS_ISNOUNMAPWAIT(as)) {
1476         if (AS_ISUNMAPWAIT(as) == 0)
1477             cv_broadcast(&as->a_cv);
1478         AS_SETUNMAPWAIT(as);
1479         AS_LOCK_EXIT(as, &as->a_lock);
1480         while (AS_ISUNMAPWAIT(as))
1481             cv_wait(&as->a_cv, &as->a_contents);
1482     } else {
1483         /*
1484          * We may have raced with
1485          * segvn_reclaim()/segspt_reclaim(). In this
1486          * case clean nounmapwait flag and retry since
1487          * softlockcnt in this segment may be already
1488          * 0. We don't drop as writer lock so our
1489          * number of retries without sleeping should
1490          * be very small. See segvn_reclaim() for
1491          * more comments.
1492          */
1493         AS_CLRNOUNMAPWAIT(as);
1494         mutex_exit(&as->a_contents);
1495         goto retry;
1496     }
1497     mutex_exit(&as->a_contents);
1498     goto top;
1499 } else if (err == IE_RETRY) {
1500     AS_LOCK_EXIT(as, &as->a_lock);
1501     goto top;
1502 } else if (err) {
1503     as_setwatch(as);
1504     AS_LOCK_EXIT(as, &as->a_lock);
1505     return (-1);
1506 }
1508     as->a_size -= ssize;
1509     if (rsize)
1510         as->a_resvsize -= rsize;
1511     raddr += ssize;
1512 }
1513     AS_LOCK_EXIT(as, &as->a_lock);
1514     return (0);
1515 }

```

unchanged_portion_omitted

```

1848 /*
1849  * Delete all segments in the address space marked with S_PURGE.
1850  * This is currently used for Sparc V9 nofault ASI segments (seg_nf.c).

```

```

1851  * These segments are deleted as a first step before calls to as_gap(), so
1852  * that they don't affect mmap() or shmat().
1853  */
1854 void
1855 as_purge(struct as *as)
1856 {
1857     struct seg *seg;
1858     struct seg *next_seg;
1859
1860     /*
1861      * the setting of NEEDSPURGE is protect by as_rangelock(), so
1862      * no need to grab a_contents mutex for this check
1863      */
1864     if ((as->a_flags & AS_NEEDSPURGE) == 0)
1865         return;
1866
1867     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
1868     next_seg = NULL;
1869     seg = AS_SEGFIRST(as);
1870     while (seg != NULL) {
1871         next_seg = AS_SEGNEXT(as, seg);
1872         if (seg->s_flags & S_PURGE)
1873             (void) segop_unmap(seg, seg->s_base, seg->s_size);
1874             SEGOP_UNMAP(seg, seg->s_base, seg->s_size);
1875         seg = next_seg;
1876     }
1877     AS_LOCK_EXIT(as, &as->a_lock);
1878
1879     mutex_enter(&as->a_contents);
1880     as->a_flags &= ~AS_NEEDSPURGE;
1881     mutex_exit(&as->a_contents);
1882 }

```

unchanged_portion_omitted

```

2143 /*
2144  * Swap the pages associated with the address space as out to
2145  * secondary storage, returning the number of bytes actually
2146  * swapped.
2147  *
2148  * The value returned is intended to correlate well with the process's
2149  * memory requirements. Its usefulness for this purpose depends on
2150  * how well the segment-level routines do at returning accurate
2151  * information.
2152  */
2153 size_t
2154 as_swapout(struct as *as)
2155 {
2156     struct seg *seg;
2157     size_t swpcnt = 0;
2158
2159     /*
2160      * Kernel-only processes have given up their address
2161      * spaces. Of course, we shouldn't be attempting to
2162      * swap out such processes in the first place...
2163      */
2164     if (as == NULL)
2165         return (0);
2166
2167     AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
2168
2169     /* Prevent XHATs from attaching */
2170     mutex_enter(&as->a_contents);
2171     AS_SETBUSY(as);
2172     mutex_exit(&as->a_contents);

```

```

2175 /*
2176  * Free all mapping resources associated with the address
2177  * space. The segment-level swapout routines capitalize
2178  * on this unmapping by scavenging pages that have become
2179  * unmapped here.
2180  */
2181 hat_swapout(as->a_hat);
2182 if (as->a_xhat != NULL)
2183     xhat_swapout_all(as);
2185 mutex_enter(&as->a_contents);
2186 AS_CLRBUSY(as);
2187 mutex_exit(&as->a_contents);
2189 /*
2190  * Call the swapout routines of all segments in the address
2191  * space to do the actual work, accumulating the amount of
2192  * space reclaimed.
2193  */
2194 for (seg = AS_SEGFIRST(as); seg != NULL; seg = AS_SEGNEXT(as, seg)) {
2195     struct seg_ops *ov = seg->s_ops;
2197     /*
2198      * We have to check to see if the seg has
2199      * an ops vector because the seg may have
2200      * been in the middle of being set up when
2201      * the process was picked for swapout.
2202      */
2203     if ((ov != NULL) && (ov->swapout != NULL))
2204         swpcnt += segop_swapout(seg);
2205     swpcnt += SEGOP_SWAPOUT(seg);
2206     AS_LOCK_EXIT(as, &as->a_lock);
2207     return (swpcnt);
2208 }
2210 /*
2211  * Determine whether data from the mappings in interval [addr, addr + size)
2212  * are in the primary memory (core) cache.
2213  */
2214 int
2215 as_incore(struct as *as, caddr_t addr,
2216           size_t size, char *vec, size_t *sizep)
2217 {
2218     struct seg *seg;
2219     size_t ssize;
2220     caddr_t raddr; /* rounded down addr */
2221     size_t rsize; /* rounded up size */
2222     size_t isize; /* iteration size */
2223     int error = 0; /* result, assume success */
2225     *sizep = 0;
2226     raddr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);
2227     rsize = (((size_t)addr + size) + PAGEOFFSET) & PAGEMASK -
2228             (size_t)raddr;
2230     if (raddr + rsize < raddr) /* check for wraparound */
2231         return (ENOMEM);
2233     AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
2234     seg = as_segat(as, raddr);
2235     if (seg == NULL) {
2236         AS_LOCK_EXIT(as, &as->a_lock);
2237         return (-1);
2238     }

```

```

2240     for (; rsize != 0; rsize -= ssize, raddr += ssize) {
2241         if (raddr >= seg->s_base + seg->s_size) {
2242             seg = AS_SEGNEXT(as, seg);
2243             if (seg == NULL || raddr != seg->s_base) {
2244                 error = -1;
2245                 break;
2246             }
2247         }
2248         if ((raddr + rsize) > (seg->s_base + seg->s_size))
2249             ssize = seg->s_base + seg->s_size - raddr;
2250         else
2251             ssize = rsize;
2252         *sizep += isize = segop_incore(seg, raddr, ssize, vec);
2253         *sizep += isize = SEGOP_INCORE(seg, raddr, ssize, vec);
2254         if (isize != ssize) {
2255             error = -1;
2256             break;
2257         }
2258         vec += btopr(ssize);
2259     }
2260     AS_LOCK_EXIT(as, &as->a_lock);
2261     return (error);
2263 static void
2264 as_segunlock(struct seg *seg, caddr_t addr, int attr,
2265             ulong_t *bitmap, size_t position, size_t npages)
2266 {
2267     caddr_t range_start;
2268     size_t pos1 = position;
2269     size_t pos2;
2270     size_t size;
2271     size_t end_pos = npages + position;
2273     while (bt_range(bitmap, &pos1, &pos2, end_pos)) {
2274         size = ptob((pos2 - pos1));
2275         range_start = (caddr_t)((uintptr_t)addr +
2276                                ptob(pos1 - position));
2278         (void) segop_lockop(seg, range_start, size, attr, MC_UNLOCK,
2279                             (void) SEGOP_LOCKOP(seg, range_start, size, attr, MC_UNLOCK,
2280                                                    (ulong_t *)NULL, (size_t)NULL);
2281         pos1 = pos2;
2282     }
2283     unchanged portion omitted
2307 /*
2308  * Cache control operations over the interval [addr, addr + size) in
2309  * address space "as".
2310  */
2311 /*ARGSUSED*/
2312 int
2313 as_ctl(struct as *as, caddr_t addr, size_t size, int func, int attr,
2314        uintptr_t arg, ulong_t *lock_map, size_t pos)
2315 {
2316     struct seg *seg; /* working segment */
2317     caddr_t raddr; /* rounded down addr */
2318     caddr_t initraddr; /* saved initial rounded down addr */
2319     size_t rsize; /* rounded up size */
2320     size_t initrsize; /* saved initial rounded up size */
2321     size_t ssize; /* size of seg */
2322     int error = 0; /* result */
2323     size_t mlock_size; /* size of bitmap */
2324     ulong_t *mlock_map; /* pointer to bitmap used */
2325     /* to represent the locked */

```

```

2326             /* pages. */
2327 retry:
2328     if (error == IE_RETRY)
2329         AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
2330     else
2331         AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
2332
2333     /*
2334     * If these are address space lock/unlock operations, loop over
2335     * all segments in the address space, as appropriate.
2336     */
2337     if (func == MC_LOCKAS) {
2338         size_t npages, idx;
2339         size_t rlen = 0;          /* rounded as length */
2340
2341         idx = pos;
2342
2343         if (arg & MCL_FUTURE) {
2344             mutex_enter(&as->a_contents);
2345             AS_SETPGLCK(as);
2346             mutex_exit(&as->a_contents);
2347         }
2348         if ((arg & MCL_CURRENT) == 0) {
2349             AS_LOCK_EXIT(as, &as->a_lock);
2350             return (0);
2351         }
2352
2353         seg = AS_SEGFIRST(as);
2354         if (seg == NULL) {
2355             AS_LOCK_EXIT(as, &as->a_lock);
2356             return (0);
2357         }
2358
2359         do {
2360             raddr = (caddr_t)((uintptr_t)seg->s_base &
2361                 (uintptr_t)PAGEMASK);
2362             rlen += (((uintptr_t)(seg->s_base + seg->s_size) +
2363                 PAGEOFFSET) & PAGEMASK) - (uintptr_t)raddr;
2364         } while ((seg = AS_SEGNEXT(as, seg)) != NULL);
2365
2366         mlock_size = BT_BITOUL(btopr(rlen));
2367         if ((mlock_map = (ulong_t *)kmem_zalloc(mlock_size *
2368             sizeof(ulong_t), KM_NOSLEEP)) == NULL) {
2369             AS_LOCK_EXIT(as, &as->a_lock);
2370             return (EAGAIN);
2371         }
2372
2373         for (seg = AS_SEGFIRST(as); seg; seg = AS_SEGNEXT(as, seg)) {
2374             error = segop_lockop(seg, seg->s_base,
2375                 error = SEGOP_LOCKOP(seg, seg->s_base,
2376                 seg->s_size, attr, MC_LOCK, mlock_map, pos);
2377             if (error != 0)
2378                 break;
2379             pos += seg_pages(seg);
2380         }
2381
2382         if (error) {
2383             for (seg = AS_SEGFIRST(as); seg != NULL;
2384                 seg = AS_SEGNEXT(as, seg)) {
2385                 raddr = (caddr_t)((uintptr_t)seg->s_base &
2386                     (uintptr_t)PAGEMASK);
2387                 npages = seg_pages(seg);
2388                 as_segunlock(seg, raddr, attr, mlock_map,
2389                     idx, npages);
2390                 idx += npages;

```

```

2391             }
2392         }
2393
2394         kmem_free(mlock_map, mlock_size * sizeof(ulong_t));
2395         AS_LOCK_EXIT(as, &as->a_lock);
2396         goto lockerr;
2397     } else if (func == MC_UNLOCKAS) {
2398         mutex_enter(&as->a_contents);
2399         AS_CLRPGGLCK(as);
2400         mutex_exit(&as->a_contents);
2401
2402         for (seg = AS_SEGFIRST(as); seg; seg = AS_SEGNEXT(as, seg)) {
2403             error = segop_lockop(seg, seg->s_base,
2404                 error = SEGOP_LOCKOP(seg, seg->s_base,
2405                 seg->s_size, attr, MC_UNLOCK, NULL, 0);
2406             if (error != 0)
2407                 break;
2408         }
2409
2410         AS_LOCK_EXIT(as, &as->a_lock);
2411         goto lockerr;
2412     }
2413
2414     /*
2415     * Normalize addresses and sizes.
2416     */
2417     initraddr = raddr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);
2418     initsize = rsize = (((size_t)(addr + size) + PAGEOFFSET) & PAGEMASK) -
        (size_t)raddr;
2419
2420     if (raddr + rsize < raddr) {          /* check for wraparound */
2421         AS_LOCK_EXIT(as, &as->a_lock);
2422         return (ENOMEM);
2423     }
2424
2425     /*
2426     * Get initial segment.
2427     */
2428     if ((seg = as_segat(as, raddr)) == NULL) {
2429         AS_LOCK_EXIT(as, &as->a_lock);
2430         return (ENOMEM);
2431     }
2432
2433     if (func == MC_LOCK) {
2434         mlock_size = BT_BITOUL(btopr(rsize));
2435         if ((mlock_map = (ulong_t *)kmem_zalloc(mlock_size *
2436             sizeof(ulong_t), KM_NOSLEEP)) == NULL) {
2437             AS_LOCK_EXIT(as, &as->a_lock);
2438             return (EAGAIN);
2439         }
2440     }
2441
2442     /*
2443     * Loop over all segments.  If a hole in the address range is
2444     * discovered, then fail.  For each segment, perform the appropriate
2445     * control operation.
2446     */
2447     while (rsize != 0) {
2448
2449         /*
2450         * Make sure there's no hole, calculate the portion
2451         * of the next segment to be operated over.
2452         */
2453         if (raddr >= seg->s_base + seg->s_size) {
2454             seg = AS_SEGNEXT(as, seg);
2455             if (seg == NULL || raddr != seg->s_base) {

```

```

2456         if (func == MC_LOCK) {
2457             as_unlockerr(as, attr, mlock_map,
2458                 initraddr, initsize - rsize);
2459             kmem_free(mlock_map,
2460                 mlock_size * sizeof (ulong_t));
2461         }
2462         AS_LOCK_EXIT(as, &as->a_lock);
2463         return (ENOMEM);
2464     }
2465 }
2466 if ((raddr + rsize) > (seg->s_base + seg->s_size))
2467     ssize = seg->s_base + seg->s_size - raddr;
2468 else
2469     ssize = rsize;
2471
2472 /*
2473  * Dispatch on specific function.
2474 */
2475 switch (func) {
2476
2477 /*
2478  * Synchronize cached data from mappings with backing
2479  * objects.
2480 */
2481 case MC_SYNC:
2482     if (error = segop_sync(seg, raddr, ssize,
2483         if (error = SEGOP_SYNC(seg, raddr, ssize,
2484             attr, (uint_t)arg)) {
2485                 AS_LOCK_EXIT(as, &as->a_lock);
2486                 return (error);
2487             }
2488         }
2489     break;
2490
2491 /*
2492  * Lock pages in memory.
2493 */
2494 case MC_LOCK:
2495     if (error = segop_lockop(seg, raddr, ssize,
2496         if (error = SEGOP_LOCKOP(seg, raddr, ssize,
2497             attr, func, mlock_map, pos)) {
2498                 as_unlockerr(as, attr, mlock_map, initraddr,
2499                     initsize - rsize + ssize);
2500                 kmem_free(mlock_map, mlock_size *
2501                     sizeof (ulong_t));
2502                 AS_LOCK_EXIT(as, &as->a_lock);
2503                 goto lockerr;
2504             }
2505         }
2506     break;
2507
2508 /*
2509  * Unlock mapped pages.
2510 */
2511 case MC_UNLOCK:
2512     (void) segop_lockop(seg, raddr, ssize, attr, func,
2513         (void) SEGOP_LOCKOP(seg, raddr, ssize, attr, func,
2514             (ulong_t *)NULL, (size_t)NULL);
2515     break;
2516
2517 /*
2518  * Store VM advise for mapped pages in segment layer.
2519 */
2520 case MC_ADVISE:
2521     error = segop_advise(seg, raddr, ssize, (uint_t)arg);
2522     error = SEGOP_ADVISE(seg, raddr, ssize, (uint_t)arg);
2523
2524 /*

```

```

2518         * Check for regular errors and special retry error
2519         */
2520     if (error) {
2521         if (error == IE_RETRY) {
2522             /*
2523              * Need to acquire writers lock, so
2524              * have to drop readers lock and start
2525              * all over again
2526              */
2527             AS_LOCK_EXIT(as, &as->a_lock);
2528             goto retry;
2529         } else if (error == IE_REATTACH) {
2530             /*
2531              * Find segment for current address
2532              * because current segment just got
2533              * split or concatenated
2534              */
2535             seg = as_segat(as, raddr);
2536             if (seg == NULL) {
2537                 AS_LOCK_EXIT(as, &as->a_lock);
2538                 return (ENOMEM);
2539             }
2540         } else {
2541             /*
2542              * Regular error
2543              */
2544             AS_LOCK_EXIT(as, &as->a_lock);
2545             return (error);
2546         }
2547     }
2548     break;
2549
2550 case MC_INHERIT_ZERO:
2551     if (seg->s_ops->inherit == NULL) {
2552         error = ENOTSUP;
2553     } else {
2554         error = segop_inherit(seg, raddr, ssize,
2555             error = SEGOP_INHERIT(seg, raddr, ssize,
2556                 SEGP_INH_ZERO);
2557     }
2558     if (error != 0) {
2559         AS_LOCK_EXIT(as, &as->a_lock);
2560         return (error);
2561     }
2562     break;
2563
2564 /*
2565  * Can't happen.
2566 */
2567 default:
2568     panic("as_ctl: bad operation %d", func);
2569     /*NOTREACHED*/
2570 }
2571
2572 rsize -= ssize;
2573 raddr += ssize;
2574
2575 if (func == MC_LOCK)
2576     kmem_free(mlock_map, mlock_size * sizeof (ulong_t));
2577 AS_LOCK_EXIT(as, &as->a_lock);
2578 return (0);
2579 lockerr:
2580
2581 /*
2582  * If the lower levels returned EDEADLK for a segment lockop,

```

```

2583      * it means that we should retry the operation. Let's wait
2584      * a bit also to let the deadlock causing condition clear.
2585      * This is part of a gross hack to work around a design flaw
2586      * in the ufs/sds logging code and should go away when the
2587      * logging code is re-designed to fix the problem. See bug
2588      * 4125102 for details of the problem.
2589      */
2590      if (error == EDEADLK) {
2591          delay(deadlk_wait);
2592          error = 0;
2593          goto retry;
2594      }
2595      return (error);
2596 }
_____ unchanged_portion_omitted _____
2617 /*
2618  * Pagelock pages from a range that spans more than 1 segment. Obtain shadow
2619  * lists from each segment and copy them to one contiguous shadow list (plist)
2620  * as expected by the caller. Save pointers to per segment shadow lists at
2621  * the tail of plist so that they can be used during as_pageunlock().
2622  */
2623 static int
2624 as_pagelock_segs(struct as *as, struct seg *seg, struct page ***ppp,
2625                caddr_t addr, size_t size, enum seg_rw rw)
2626 {
2627     caddr_t sv_addr = addr;
2628     size_t sv_size = size;
2629     struct seg *sv_seg = seg;
2630     ulong_t segcnt = 1;
2631     ulong_t cnt;
2632     size_t ssize;
2633     pgcnt_t npages = btop(size);
2634     page_t **plist;
2635     page_t **pl;
2636     int error;
2637     caddr_t eaddr;
2638     faultcode_t fault_err = 0;
2639     pgcnt_t pl_off;
2640     extern struct seg_ops segspt_shmops;

2642     ASSERT(AS_LOCK_HELD(as, &as->a_lock));
2643     ASSERT(seg != NULL);
2644     ASSERT(addr >= seg->s_base && addr < seg->s_base + seg->s_size);
2645     ASSERT(addr + size > seg->s_base + seg->s_size);
2646     ASSERT(IS_P2ALIGNED(size, PAGESIZE));
2647     ASSERT(IS_P2ALIGNED(addr, PAGESIZE));

2649     /*
2650     * Count the number of segments covered by the range we are about to
2651     * lock. The segment count is used to size the shadow list we return
2652     * back to the caller.
2653     */
2654     for (; size != 0; size -= ssize, addr += ssize) {
2655         if (addr >= seg->s_base + seg->s_size) {

2657             seg = AS_SEGNEXT(as, seg);
2658             if (seg == NULL || addr != seg->s_base) {
2659                 AS_LOCK_EXIT(as, &as->a_lock);
2660                 return (EFAULT);
2661             }
2662             /*
2663             * Do a quick check if subsequent segments
2664             * will most likely support pagelock.
2665             */
2666             if (seg->s_ops == &segvn_ops) {

```

```

2667         vnode_t *vp;

2669         if (segop_getvp(seg, addr, &vp) != 0 ||
2669         if (SEGOP_GETVP(seg, addr, &vp) != 0 ||
2670             vp != NULL) {
2671             AS_LOCK_EXIT(as, &as->a_lock);
2672             goto slow;
2673         }
2674         } else if (seg->s_ops != &segspt_shmops) {
2675             AS_LOCK_EXIT(as, &as->a_lock);
2676             goto slow;
2677         }
2678         segcnt++;
2679     }
2680     if (addr + size > seg->s_base + seg->s_size) {
2681         ssize = seg->s_base + seg->s_size - addr;
2682     } else {
2683         ssize = size;
2684     }
2685     }
2686     ASSERT(segcnt > 1);

2688     plist = kmem_zalloc((npages + segcnt) * sizeof (page_t *), KM_SLEEP);

2690     addr = sv_addr;
2691     size = sv_size;
2692     seg = sv_seg;

2694     for (cnt = 0, pl_off = 0; size != 0; size -= ssize, addr += ssize) {
2695         if (addr >= seg->s_base + seg->s_size) {
2696             seg = AS_SEGNEXT(as, seg);
2697             ASSERT(seg != NULL && addr == seg->s_base);
2698             cnt++;
2699             ASSERT(cnt < segcnt);
2700         }
2701         if (addr + size > seg->s_base + seg->s_size) {
2702             ssize = seg->s_base + seg->s_size - addr;
2703         } else {
2704             ssize = size;
2705         }
2706         pl = &plist[npages + cnt];
2707         error = segop_pagelock(seg, addr, ssize, (page_t ***)pl,
2707         error = SEGOP_PAGELOCK(seg, addr, ssize, (page_t ***)pl,
2708             L_PAGELOCK, rw);
2709         if (error) {
2710             break;
2711         }
2712         ASSERT(plist[npages + cnt] != NULL);
2713         ASSERT(pl_off + btop(ssize) <= npages);
2714         bcopy(plist[npages + cnt], &plist[pl_off],
2715             btop(ssize) * sizeof (page_t *));
2716         pl_off += btop(ssize);
2717     }

2719     if (size == 0) {
2720         AS_LOCK_EXIT(as, &as->a_lock);
2721         ASSERT(cnt == segcnt - 1);
2722         *ppp = plist;
2723         return (0);
2724     }

2726     /*
2727     * one of pagelock calls failed. The error type is in error variable.
2728     * Unlock what we've locked so far and retry with F_SOFTLOCK if error
2729     * type is either EFAULT or ENOTSUP. Otherwise just return the error
2730     * back to the caller.

```

```

2731     */
2732
2733     eaddr = addr;
2734     seg = sv_seg;
2735
2736     for (cnt = 0, addr = sv_addr; addr < eaddr; addr += ssize) {
2737         if (addr >= seg->s_base + seg->s_size) {
2738             seg = AS_SEGNEXT(as, seg);
2739             ASSERT(seg != NULL && addr == seg->s_base);
2740             cnt++;
2741             ASSERT(cnt < segcnt);
2742         }
2743         if (eaddr > seg->s_base + seg->s_size) {
2744             ssize = seg->s_base + seg->s_size - addr;
2745         } else {
2746             ssize = eaddr - addr;
2747         }
2748         pl = &plist[npages + cnt];
2749         ASSERT(*pl != NULL);
2750         (void) segop_pagelock(seg, addr, ssize, (page_t ***)pl,
2751             (void) SEGOP_PAGELOCK(seg, addr, ssize, (page_t ***)pl,
2752             L_PAGEUNLOCK, rw));
2753     }
2754     AS_LOCK_EXIT(as, &as->a_lock);
2755
2756     kmem_free(plist, (npages + segcnt) * sizeof (page_t *));
2757
2758     if (error != ENOTSUP && error != EFAULT) {
2759         return (error);
2760     }
2761
2762 slow:
2763     /*
2764     * If we are here because pagelock failed due to the need to cow fault
2765     * in the pages we want to lock F_SOFTLOCK will do this job and in
2766     * next as_pagelock() call for this address range pagelock will
2767     * hopefully succeed.
2768     */
2769     fault_err = as_fault(as->a_hat, as, sv_addr, sv_size, F_SOFTLOCK, rw);
2770     if (fault_err != 0) {
2771         return (fc_decode(fault_err));
2772     }
2773     *ppp = NULL;
2774
2775     return (0);
2776 }
2777
2778 /*
2779 * lock pages in a given address space. Return shadow list. If
2780 * the list is NULL, the MMU mapping is also locked.
2781 */
2782 int
2783 as_pagelock(struct as *as, struct page **ppp, caddr_t addr,
2784     size_t size, enum seg_rw rw)
2785 {
2786     size_t rsize;
2787     caddr_t raddr;
2788     faultcode_t fault_err;
2789     struct seg *seg;
2790     int err;
2791
2792     TRACE_2(TR_FAC_PHYSIO, TR_PHYSIO_AS_LOCK_START,
2793         "as_pagelock_start: addr %p size %ld", addr, size);
2794
2795     raddr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);

```

```

2796     rsize = (((size_t)(addr + size) + PAGEOFFSET) & PAGEMASK) -
2797         (size_t)raddr;
2798
2799     /*
2800     * if the request crosses two segments let
2801     * as_fault handle it.
2802     */
2803     AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
2804
2805     seg = as_segat(as, raddr);
2806     if (seg == NULL) {
2807         AS_LOCK_EXIT(as, &as->a_lock);
2808         return (EFAULT);
2809     }
2810     ASSERT(raddr >= seg->s_base && raddr < seg->s_base + seg->s_size);
2811     if (raddr + rsize > seg->s_base + seg->s_size) {
2812         return (as_pagelock_segs(as, seg, ppp, raddr, rsize, rw));
2813     }
2814     if (raddr + rsize <= raddr) {
2815         AS_LOCK_EXIT(as, &as->a_lock);
2816         return (EFAULT);
2817     }
2818
2819     TRACE_2(TR_FAC_PHYSIO, TR_PHYSIO_SEG_LOCK_START,
2820         "seg_lock_l_start: raddr %p rsize %ld", raddr, rsize);
2821
2822     /*
2823     * try to lock pages and pass back shadow list
2824     */
2825     err = segop_pagelock(seg, raddr, rsize, ppp, L_PAGELOCK, rw);
2826     err = SEGOP_PAGELOCK(seg, raddr, rsize, ppp, L_PAGELOCK, rw);
2827
2828     TRACE_0(TR_FAC_PHYSIO, TR_PHYSIO_SEG_LOCK_END, "seg_lock_l_end");
2829
2830     AS_LOCK_EXIT(as, &as->a_lock);
2831
2832     if (err == 0 || (err != ENOTSUP && err != EFAULT)) {
2833         return (err);
2834     }
2835
2836     /*
2837     * Use F_SOFTLOCK to lock the pages because pagelock failed either due
2838     * to no pagelock support for this segment or pages need to be cow
2839     * faulted in. If fault is needed F_SOFTLOCK will do this job for
2840     * this as_pagelock() call and in the next as_pagelock() call for the
2841     * same address range pagelock call will hopefully succeed.
2842     */
2843     fault_err = as_fault(as->a_hat, as, addr, size, F_SOFTLOCK, rw);
2844     if (fault_err != 0) {
2845         return (fc_decode(fault_err));
2846     }
2847     *ppp = NULL;
2848
2849     TRACE_0(TR_FAC_PHYSIO, TR_PHYSIO_AS_LOCK_END, "as_pagelock_end");
2850     return (0);
2851 }
2852
2853 /*
2854 * unlock pages locked by as_pagelock_segs(). Retrieve per segment shadow
2855 * lists from the end of plist and call pageunlock interface for each segment.
2856 * Drop as lock and free plist.
2857 */
2858 static void
2859 as_pageunlock_segs(struct as *as, struct seg *seg, caddr_t addr, size_t size,
2860     struct page **plist, enum seg_rw rw)

```

```

2861     ulong_t cnt;
2862     caddr_t eaddr = addr + size;
2863     pgcnt_t npages = btop(size);
2864     size_t ssize;
2865     page_t **pl;

2867     ASSERT(AS_LOCK_HELD(as, &as->a_lock));
2868     ASSERT(seg != NULL);
2869     ASSERT(addr >= seg->s_base && addr < seg->s_base + seg->s_size);
2870     ASSERT(addr + size > seg->s_base + seg->s_size);
2871     ASSERT(IS_P2ALIGNED(size, PAGE_SIZE));
2872     ASSERT(IS_P2ALIGNED(addr, PAGE_SIZE));
2873     ASSERT(plist != NULL);

2875     for (cnt = 0; addr < eaddr; addr += ssize) {
2876         if (addr >= seg->s_base + seg->s_size) {
2877             seg = AS_SEGNEXT(as, seg);
2878             ASSERT(seg != NULL && addr == seg->s_base);
2879             cnt++;
2880         }
2881         if (eaddr > seg->s_base + seg->s_size) {
2882             ssize = seg->s_base + seg->s_size - addr;
2883         } else {
2884             ssize = eaddr - addr;
2885         }
2886         pl = &plist[npages + cnt];
2887         ASSERT(*pl != NULL);
2888         (void) segop_pagelock(seg, addr, ssize, (page_t ***)pl,
2888             (void) SEGOP_PAGELOCK(seg, addr, ssize, (page_t ***)pl,
2889                 L_PAGEUNLOCK, rw));
2890     }
2891     ASSERT(cnt > 0);
2892     AS_LOCK_EXIT(as, &as->a_lock);

2894     cnt++;
2895     kmem_free(plist, (npages + cnt) * sizeof (page_t *));
2896 }

2898 /*
2899  * unlock pages in a given address range
2900  */
2901 void
2902 as_pageunlock(struct as *as, struct page **pp, caddr_t addr, size_t size,
2903     enum seg_rw rw)
2904 {
2905     struct seg *seg;
2906     size_t rsize;
2907     caddr_t raddr;

2909     TRACE_2(TR_FAC_PHYSIO, TR_PHYSIO_AS_UNLOCK_START,
2910         "as_pageunlock_start: addr %p size %ld", addr, size);

2912     /*
2913      * if the shadow list is NULL, as_pagelock was
2914      * falling back to as_fault
2915      */
2916     if (pp == NULL) {
2917         (void) as_fault(as->a_hat, as, addr, size, F_SOFTUNLOCK, rw);
2918         return;
2919     }

2921     raddr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);
2922     rsize = (((size_t)(addr + size) + PAGEOFFSET) & PAGEMASK) -
2923         (size_t)raddr;

2925     AS_LOCK_ENTER(as, &as->a_lock, RW_READER);

```

```

2926     seg = as_segat(as, raddr);
2927     ASSERT(seg != NULL);

2929     TRACE_2(TR_FAC_PHYSIO, TR_PHYSIO_SEG_UNLOCK_START,
2930         "seg_unlock_start: raddr %p rsize %ld", raddr, rsize);

2932     ASSERT(raddr >= seg->s_base && raddr < seg->s_base + seg->s_size);
2933     if (raddr + rsize <= seg->s_base + seg->s_size) {
2934         (void) segop_pagelock(seg, raddr, rsize, &pp, L_PAGEUNLOCK, rw);
2935     } else {
2936         (void) SEGOP_PAGELOCK(seg, raddr, rsize, &pp, L_PAGEUNLOCK, rw);
2937     } else {
2938         as_pageunlock_segs(as, seg, raddr, rsize, pp, rw);
2939         return;
2940     }
2941     AS_LOCK_EXIT(as, &as->a_lock);
2942     TRACE_0(TR_FAC_PHYSIO, TR_PHYSIO_AS_UNLOCK_END, "as_pageunlock_end");

2943 int
2944 as_setpagesize(struct as *as, caddr_t addr, size_t size, uint_t szc,
2945     boolean_t wait)
2946 {
2947     struct seg *seg;
2948     size_t ssize;
2949     caddr_t raddr;
2950     size_t rsize;
2951     int error = 0;
2952     size_t pgsz = page_get_pagesize(szc);

2954     setpgsz_top:
2955     if (!IS_P2ALIGNED(addr, pgsz) || !IS_P2ALIGNED(size, pgsz)) {
2956         return (EINVAL);
2957     }

2959     raddr = addr;
2960     rsize = size;

2962     if (raddr + rsize < raddr)
2963         return (ENOMEM); /* check for wraparound */

2965     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
2966     as_clearwatchprot(as, raddr, rsize);
2967     seg = as_segat(as, raddr);
2968     if (seg == NULL) {
2969         as_setwatch(as);
2970         AS_LOCK_EXIT(as, &as->a_lock);
2971         return (ENOMEM);
2972     }

2974     for (; rsize != 0; rsize -= ssize, raddr += ssize) {
2975         if (raddr >= seg->s_base + seg->s_size) {
2976             seg = AS_SEGNEXT(as, seg);
2977             if (seg == NULL || raddr != seg->s_base) {
2978                 error = ENOMEM;
2979                 break;
2980             }
2981         }
2982         if ((raddr + rsize) > (seg->s_base + seg->s_size)) {
2983             ssize = seg->s_base + seg->s_size - raddr;
2984         } else {
2985             ssize = rsize;
2986         }

2988     retry:
2989     error = segop_setpagesize(seg, raddr, ssize, szc);
2990     error = SEGOP_SETPAGESIZE(seg, raddr, ssize, szc);

```

```

2991     if (error == IE_NOMEM) {
2992         error = EAGAIN;
2993         break;
2994     }
2995
2996     if (error == IE_RETRY) {
2997         AS_LOCK_EXIT(as, &as->a_lock);
2998         goto setpgsz_top;
2999     }
3000
3001     if (error == ENOTSUP) {
3002         error = EINVAL;
3003         break;
3004     }
3005
3006     if (wait && (error == EAGAIN)) {
3007         /*
3008          * Memory is currently locked. It must be unlocked
3009          * before this operation can succeed through a retry.
3010          * The possible reasons for locked memory and
3011          * corresponding strategies for unlocking are:
3012          * (1) Normal I/O
3013          *     wait for a signal that the I/O operation
3014          *     has completed and the memory is unlocked.
3015          * (2) Asynchronous I/O
3016          *     The aio subsystem does not unlock pages when
3017          *     the I/O is completed. Those pages are unlocked
3018          *     when the application calls aiowait/aioerror.
3019          *     So, to prevent blocking forever, cv_broadcast()
3020          *     is done to wake up aio_cleanup_thread.
3021          *     Subsequently, segvn_reclaim will be called, and
3022          *     that will do AS_CLRUNMAPWAIT() and wake us up.
3023          * (3) Long term page locking:
3024          *     This is not relevant for as_setpagesize()
3025          *     because we cannot change the page size for
3026          *     driver memory. The attempt to do so will
3027          *     fail with a different error than EAGAIN so
3028          *     there's no need to trigger as callbacks like
3029          *     as_unmap, as_setprot or as_free would do.
3030          */
3031         mutex_enter(&as->a_contents);
3032         if (!AS_ISNOUNMAPWAIT(as)) {
3033             if (AS_ISUNMAPWAIT(as) == 0) {
3034                 cv_broadcast(&as->a_cv);
3035             }
3036             AS_SETUNMAPWAIT(as);
3037             AS_LOCK_EXIT(as, &as->a_lock);
3038             while (AS_ISUNMAPWAIT(as)) {
3039                 cv_wait(&as->a_cv, &as->a_contents);
3040             }
3041         } else {
3042             /*
3043              * We may have raced with
3044              * segvn_reclaim()/segspt_reclaim(). In this
3045              * case clean nounmapwait flag and retry since
3046              * softlockt in this segment may be already
3047              * 0. We don't drop as writer lock so our
3048              * number of retries without sleeping should
3049              * be very small. See segvn_reclaim() for
3050              * more comments.
3051              */
3052             AS_CLRNOUNMAPWAIT(as);
3053             mutex_exit(&as->a_contents);
3054             goto retry;
3055         }

```

```

3056         mutex_exit(&as->a_contents);
3057         goto setpgsz_top;
3058     } else if (error != 0) {
3059         break;
3060     }
3061 }
3062 as_setwatch(as);
3063 AS_LOCK_EXIT(as, &as->a_lock);
3064 return (error);
3065 }
3066
3067 /*
3068  * as_iset3_default_lpsize() just calls segop_setpagesize() on all segments
3069  * as_iset3_default_lpsize() just calls SEGOP_SETPAGESIZE() on all segments
3070  * in its chunk where s_szc is less than the szc we want to set.
3071  */
3072 static int
3073 as_iset3_default_lpsize(struct as *as, caddr_t raddr, size_t rsize, uint_t szc,
3074 int *retry)
3075 {
3076     struct seg *seg;
3077     size_t ssize;
3078     int error;
3079
3080     ASSERT(AS_WRITE_HELD(as, &as->a_lock));
3081
3082     seg = as_segat(as, raddr);
3083     if (seg == NULL) {
3084         panic("as_iset3_default_lpsize: no seg");
3085     }
3086
3087     for (; rsize != 0; rsize -= ssize, raddr += ssize) {
3088         if (raddr >= seg->s_base + seg->s_size) {
3089             seg = AS_SEGNEXT(as, seg);
3090             if (seg == NULL || raddr != seg->s_base) {
3091                 panic("as_iset3_default_lpsize: as changed");
3092             }
3093         }
3094         if ((raddr + rsize) > (seg->s_base + seg->s_size)) {
3095             ssize = seg->s_base + seg->s_size - raddr;
3096         } else {
3097             ssize = rsize;
3098         }
3099
3100         if (szc > seg->s_szc) {
3101             error = segop_setpagesize(seg, raddr, ssize, szc);
3102             error = SEGOP_SETPAGESIZE(seg, raddr, ssize, szc);
3103             /* Only retry on EINVAL segments that have no vnode. */
3104             if (error == EINVAL) {
3105                 vnode_t *vp = NULL;
3106                 if ((segop_gettype(seg, raddr) & MAP_SHARED) &&
3107                     (segop_getvp(seg, raddr, &vp) != 0 ||
3108                     if ((SEGOP_GETTYPE(seg, raddr) & MAP_SHARED) &&
3109                         (SEGOP_GETVP(seg, raddr, &vp) != 0 ||
3110                             vp == NULL)) {
3111                     *retry = 1;
3112                 } else {
3113                     *retry = 0;
3114                 }
3115             }
3116         }
3117     }
3118     if (error) {
3119         return (error);
3120     }
3121 }
3122 return (0);

```



```

3118 }
      unchanged_portion_omitted
3305 /*
3306 * Set the default large page size for the range. Called via memcntl with
3307 * page size set to 0. as_set_default_lpsize breaks the range down into
3308 * chunks with the same type/flags, ignores-non segvn segments, and passes
3309 * each chunk to as_iset_default_lpsize().
3310 */
3311 int
3312 as_set_default_lpsize(struct as *as, caddr_t addr, size_t size)
3313 {
3314     struct seg *seg;
3315     caddr_t raddr;
3316     size_t rsize;
3317     size_t ssize;
3318     int rtype, rflags;
3319     int stype, sflags;
3320     int error;
3321     caddr_t setaddr;
3322     size_t setsize;
3323     int segvn;

3325     if (size == 0)
3326         return (0);

3328     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
3329 again:
3330     error = 0;

3332     raddr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);
3333     rsize = (((size_t)(addr + size) + PAGEOFFSET) & PAGEMASK) -
3334         (size_t)raddr;

3336     if (raddr + rsize < raddr) { /* check for wraparound */
3337         AS_LOCK_EXIT(as, &as->a_lock);
3338         return (ENOMEM);
3339     }
3340     as_clearwatchprot(as, raddr, rsize);
3341     seg = as_segat(as, raddr);
3342     if (seg == NULL) {
3343         as_setwatch(as);
3344         AS_LOCK_EXIT(as, &as->a_lock);
3345         return (ENOMEM);
3346     }
3347     if (seg->s_ops == &segvn_ops) {
3348         rtype = segop_gettype(seg, addr);
3349         rtype = SEGOP_GETTYPE(seg, addr);
3350         rflags = rtype & (MAP_TEXT | MAP_INITDATA);
3351         rtype = rtype & (MAP_SHARED | MAP_PRIVATE);
3352         segvn = 1;
3353     } else {
3354         segvn = 0;
3355     }
3356     setaddr = raddr;
3357     setsize = 0;

3358     for (; rsize != 0; rsize -= ssize, raddr += ssize, setsize += ssize) {
3359         if (raddr >= (seg->s_base + seg->s_size)) {
3360             seg = AS_SEGNEXT(as, seg);
3361             if (seg == NULL || raddr != seg->s_base) {
3362                 error = ENOMEM;
3363                 break;
3364             }
3365             if (seg->s_ops == &segvn_ops) {
3366                 stype = segop_gettype(seg, raddr);

```

```

3366         stype = SEGOP_GETTYPE(seg, raddr);
3367         sflags = stype & (MAP_TEXT | MAP_INITDATA);
3368         stype &= (MAP_SHARED | MAP_PRIVATE);
3369         if (segvn && (rflags != sflags ||
3370             rtype != stype)) {
3371             /*
3372              * The next segment is also segvn but
3373              * has different flags and/or type.
3374              */
3375             ASSERT(setsize != 0);
3376             error = as_iset_default_lpsize(as,
3377                 setaddr, setsize, rflags, rtype);
3378             if (error) {
3379                 break;
3380             }
3381             rflags = sflags;
3382             rtype = stype;
3383             setaddr = raddr;
3384             setsize = 0;
3385         } else if (!segvn) {
3386             rflags = sflags;
3387             rtype = stype;
3388             setaddr = raddr;
3389             setsize = 0;
3390             segvn = 1;
3391         }
3392     } else if (segvn) {
3393         /* The next segment is not segvn. */
3394         ASSERT(setsize != 0);
3395         error = as_iset_default_lpsize(as,
3396             setaddr, setsize, rflags, rtype);
3397         if (error) {
3398             break;
3399         }
3400         segvn = 0;
3401     }
3402 }
3403 if ((raddr + rsize) > (seg->s_base + seg->s_size)) {
3404     ssize = seg->s_base + seg->s_size - raddr;
3405 } else {
3406     ssize = rsize;
3407 }
3408 }
3409 if (error == 0 && segvn) {
3410     /* The last chunk when rsize == 0. */
3411     ASSERT(setsize != 0);
3412     error = as_iset_default_lpsize(as, setaddr, setsize,
3413         rflags, rtype);
3414 }

3416 if (error == IE_RETRY) {
3417     goto again;
3418 } else if (error == IE_NOMEM) {
3419     error = EAGAIN;
3420 } else if (error == ENOTSUP) {
3421     error = EINVAL;
3422 } else if (error == EAGAIN) {
3423     mutex_enter(&as->a_contents);
3424     if (!AS_ISNOUNMAPWAIT(as)) {
3425         if (AS_ISUNMAPWAIT(as) == 0) {
3426             cv_broadcast(&as->a_cv);
3427         }
3428         AS_SETUNMAPWAIT(as);
3429         AS_LOCK_EXIT(as, &as->a_lock);
3430         while (AS_ISUNMAPWAIT(as)) {
3431             cv_wait(&as->a_cv, &as->a_contents);

```

```

3432     }
3433     mutex_exit(&as->a_contents);
3434     AS_LOCK_ENTER(as, &as->a_lock, RW_WRITER);
3435 } else {
3436     /*
3437     * We may have raced with
3438     * segvn_reclaim()/segspt_reclaim(). In this case
3439     * clean nounmapwait flag and retry since softlocknt
3440     * in this segment may be already 0. We don't drop as
3441     * writer lock so our number of retries without
3442     * sleeping should be very small. See segvn_reclaim()
3443     * for more comments.
3444     */
3445     AS_CLRNOUNMAPWAIT(as);
3446     mutex_exit(&as->a_contents);
3447 }
3448 goto again;
3449 }

3451 as_setwatch(as);
3452 AS_LOCK_EXIT(as, &as->a_lock);
3453 return (error);
3454 }

3456 /*
3457 * Setup all of the uninitialized watched pages that we can.
3458 */
3459 void
3460 as_setwatch(struct as *as)
3461 {
3462     struct watched_page *pwp;
3463     struct seg *seg;
3464     caddr_t vaddr;
3465     uint_t prot;
3466     int err, retrycnt;

3468     if (avl_numnodes(&as->a_wpage) == 0)
3469         return;

3471     ASSERT(AS_WRITE_HELD(as, &as->a_lock));

3473     for (pwp = avl_first(&as->a_wpage); pwp != NULL;
3474          pwp = AVL_NEXT(&as->a_wpage, pwp)) {
3475         retrycnt = 0;
3476     retry:
3477         vaddr = pwp->wp_vaddr;
3478         if (pwp->wp_oprot != 0 || /* already set up */
3479             (seg = as_segat(as, vaddr)) == NULL ||
3480             segop_getprot(seg, vaddr, 0, &prot) != 0)
3481             SEGOP_GETPROT(seg, vaddr, 0, &prot) != 0)
3482             continue;

3483         pwp->wp_oprot = prot;
3484         if (pwp->wp_read)
3485             prot &= ~(PROT_READ|PROT_WRITE|PROT_EXEC);
3486         if (pwp->wp_write)
3487             prot &= ~PROT_WRITE;
3488         if (pwp->wp_exec)
3489             prot &= ~(PROT_READ|PROT_WRITE|PROT_EXEC);
3490         if (!(pwp->wp_flags & WP_NOWATCH) && prot != pwp->wp_oprot) {
3491             err = segop_setprot(seg, vaddr, PAGE_SIZE, prot);
3492             err = SEGOP_SETPROT(seg, vaddr, PAGE_SIZE, prot);
3493             if (err == IE_RETRY) {
3494                 pwp->wp_oprot = 0;
3495                 ASSERT(retrycnt == 0);
3496                 retrycnt++;

```

```

3496         goto retry;
3497     }
3498     }
3499     pwp->wp_prot = prot;
3500 }
3501 }

3503 /*
3504 * Clear all of the watched pages in the address space.
3505 */
3506 void
3507 as_clearwatch(struct as *as)
3508 {
3509     struct watched_page *pwp;
3510     struct seg *seg;
3511     caddr_t vaddr;
3512     uint_t prot;
3513     int err, retrycnt;

3515     if (avl_numnodes(&as->a_wpage) == 0)
3516         return;

3518     ASSERT(AS_WRITE_HELD(as, &as->a_lock));

3520     for (pwp = avl_first(&as->a_wpage); pwp != NULL;
3521          pwp = AVL_NEXT(&as->a_wpage, pwp)) {
3522         retrycnt = 0;
3523     retry:
3524         vaddr = pwp->wp_vaddr;
3525         if (pwp->wp_oprot == 0 || /* not set up */
3526             (seg = as_segat(as, vaddr)) == NULL)
3527             continue;

3529         if ((prot = pwp->wp_oprot) != pwp->wp_prot) {
3530             err = segop_setprot(seg, vaddr, PAGE_SIZE, prot);
3531             err = SEGOP_SETPROT(seg, vaddr, PAGE_SIZE, prot);
3532             if (err == IE_RETRY) {
3533                 ASSERT(retrycnt == 0);
3534                 retrycnt++;
3535                 goto retry;
3536             }
3537             pwp->wp_oprot = 0;
3538             pwp->wp_prot = 0;
3539         }
3540     }

3542 /*
3543 * Force a new setup for all the watched pages in the range.
3544 */
3545 static void
3546 as_setwatchprot(struct as *as, caddr_t addr, size_t size, uint_t prot)
3547 {
3548     struct watched_page *pwp;
3549     struct watched_page tpw;
3550     caddr_t eaddr = addr + size;
3551     caddr_t vaddr;
3552     struct seg *seg;
3553     int err, retrycnt;
3554     uint_t wprot;
3555     avl_index_t where;

3557     if (avl_numnodes(&as->a_wpage) == 0)
3558         return;

3560     ASSERT(AS_WRITE_HELD(as, &as->a_lock));

```

```

3562     tpw.wp_vaddr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);
3563     if ((pwp = avl_find(&as->a_wpage, &tpw, &where)) == NULL)
3564         pwp = avl_nearest(&as->a_wpage, where, AVL_AFTER);

3566     while (pwp != NULL && pwp->wp_vaddr < eaddr) {
3567         retrycnt = 0;
3568         vaddr = pwp->wp_vaddr;

3570         wprot = prot;
3571         if (pwp->wp_read)
3572             wprot &= ~(PROT_READ|PROT_WRITE|PROT_EXEC);
3573         if (pwp->wp_write)
3574             wprot &= ~PROT_WRITE;
3575         if (pwp->wp_exec)
3576             wprot &= ~(PROT_READ|PROT_WRITE|PROT_EXEC);
3577         if (!(pwp->wp_flags & WP_NOWATCH) && wprot != pwp->wp_oprot) {
3578             retry:
3579                 seg = as_segat(as, vaddr);
3580                 if (seg == NULL) {
3581                     panic("as_setwatchprot: no seg");
3582                     /*NOTREACHED*/
3583                 }
3584                 err = segop_setprot(seg, vaddr, PAGE_SIZE, wprot);
3585                 err = SEGOP_SETPROT(seg, vaddr, PAGE_SIZE, wprot);
3586                 if (err == IE_RETRY) {
3587                     ASSERT(retrycnt == 0);
3588                     retrycnt++;
3589                     goto retry;
3590                 }
3591                 pwp->wp_oprot = prot;
3592                 pwp->wp_prot = wprot;
3594             }
3595         }
3596     }

3598 /*
3599  * Clear all of the watched pages in the range.
3600  */
3601 static void
3602 as_clearwatchprot(struct as *as, caddr_t addr, size_t size)
3603 {
3604     caddr_t eaddr = addr + size;
3605     struct watched_page *pwp;
3606     struct watched_page tpw;
3607     uint_t prot;
3608     struct seg *seg;
3609     int err, retrycnt;
3610     avl_index_t where;

3612     if (avl_numnodes(&as->a_wpage) == 0)
3613         return;

3615     tpw.wp_vaddr = (caddr_t)((uintptr_t)addr & (uintptr_t)PAGEMASK);
3616     if ((pwp = avl_find(&as->a_wpage, &tpw, &where)) == NULL)
3617         pwp = avl_nearest(&as->a_wpage, where, AVL_AFTER);

3619     ASSERT(AS_WRITE_HELD(as, &as->a_lock));

3621     while (pwp != NULL && pwp->wp_vaddr < eaddr) {

3623         if ((prot = pwp->wp_oprot) != 0) {
3624             retrycnt = 0;

```

```

3626         if (prot != pwp->wp_prot) {
3627             retry:
3628                 seg = as_segat(as, pwp->wp_vaddr);
3629                 if (seg == NULL)
3630                     continue;
3631                 err = segop_setprot(seg, pwp->wp_vaddr,
3632                                     err = SEGOP_SETPROT(seg, pwp->wp_vaddr,
3633                                                         PAGE_SIZE, prot);
3634                 if (err == IE_RETRY) {
3635                     ASSERT(retrycnt == 0);
3636                     retrycnt++;
3637                     goto retry;
3638                 }
3639             }
3640             pwp->wp_oprot = 0;
3641             pwp->wp_prot = 0;
3642         }

3644         pwp = AVL_NEXT(&as->a_wpage, pwp);
3645     }
3646 }
    unchanged portion omitted

3665 /*
3666  * return memory object ID
3667  */
3668 int
3669 as_getmemid(struct as *as, caddr_t addr, memid_t *memidp)
3670 {
3671     struct seg *seg;
3672     int sts;

3674     AS_LOCK_ENTER(as, &as->a_lock, RW_READER);
3675     seg = as_segat(as, addr);
3676     if (seg == NULL) {
3677         AS_LOCK_EXIT(as, &as->a_lock);
3678         return (EFAULT);
3679     }
3680     /*
3681      * catch old drivers which may not support getmemid
3682      */
3683     if (seg->s_ops->getmemid == NULL) {
3684         AS_LOCK_EXIT(as, &as->a_lock);
3685         return (ENODEV);
3686     }

3688     sts = segop_getmemid(seg, addr, memidp);
3689     sts = SEGOP_GETMEMID(seg, addr, memidp);

3690     AS_LOCK_EXIT(as, &as->a_lock);
3691     return (sts);
3692 }
    unchanged portion omitted

```

new/usr/src/uts/common/vm/vm_pvn.c

1

```
*****
31976 Tue Nov 24 09:34:49 2015
new/usr/src/uts/common/vm/vm_pvn.c
patch lower-case-segops
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 1986, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright 2015 Nexenta Systems, Inc. All rights reserved.
24 */

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

29 /*
30 * University Copyright- Copyright (c) 1982, 1986, 1988
31 * The Regents of the University of California
32 * All Rights Reserved
33 *
34 * University Acknowledgment- Portions of this document are derived from
35 * software developed by the University of California, Berkeley, and its
36 * contributors.
37 */

39 /*
40 * VM - paged vnode.
41 *
42 * This file supplies vm support for the vnode operations that deal with pages.
43 */
44 #include <sys/types.h>
45 #include <sys/t_lock.h>
46 #include <sys/param.h>
47 #include <sys/sysmacros.h>
48 #include <sys/system.h>
49 #include <sys/time.h>
50 #include <sys/buf.h>
51 #include <sys/vnode.h>
52 #include <sys/uio.h>
53 #include <sys/vmsystem.h>
54 #include <sys/mman.h>
55 #include <sys/vfs.h>
56 #include <sys/cred.h>
57 #include <sys/user.h>
58 #include <sys/kmem.h>
59 #include <sys/cmn_err.h>
60 #include <sys/debug.h>
61 #include <sys/cpuvar.h>
```

new/usr/src/uts/common/vm/vm_pvn.c

2

```
62 #include <sys/vtrace.h>
63 #include <sys/tnf_probe.h>

65 #include <vm/hat.h>
66 #include <vm/as.h>
67 #include <vm/seg.h>
68 #include <vm/rm.h>
69 #include <vm/pvn.h>
70 #include <vm/page.h>
71 #include <vm/seg_map.h>
72 #include <vm/seg_kmem.h>
73 #include <sys/fs/swapnode.h>

75 int pvn_nofodklust = 0;
76 int pvn_write_noklust = 0;

78 uint_t pvn_vmmodsort_supported = 0;      /* set if HAT supports VMODSORT */
79 uint_t pvn_vmmodsort_disable = 0;      /* set in /etc/system to disable HAT */
80                                          /* support for vmmodsort for testing */

82 static struct kmem_cache *marker_cache = NULL;

84 /*
85 * Find the largest contiguous block which contains 'addr' for file offset
86 * 'offset' in it while living within the file system block sizes ('vp_off'
87 * and 'vp_len') and the address space limits for which no pages currently
88 * exist and which map to consecutive file offsets.
89 */
90 page_t *
91 pvn_read_kluster(
92     struct vnode *vp,
93     u_offset_t off,
94     struct seg *seg,
95     caddr_t addr,
96     u_offset_t *offp,                          /* return values */
97     size_t *lenp,                              /* return values */
98     u_offset_t vp_off,
99     size_t vp_len,
100     int isra)
101 {
102     ssize_t deltaf, deltab;
103     page_t *pp;
104     page_t *plist = NULL;
105     spgcnt_t pagesavail;
106     u_offset_t vp_end;

108     ASSERT(off >= vp_off && off < vp_off + vp_len);

110     /*
111     * We only want to do klustering/read ahead if there
112     * is more than minfree pages currently available.
113     */
114     pagesavail = freemem - minfree;

116     if (pagesavail <= 0)
117         if (isra)
118             return ((page_t *)NULL);      /* ra case - give up */
119         else
120             pagesavail = 1;                /* must return a page */

122     /* We calculate in pages instead of bytes due to 32-bit overflows */
123     if (pagesavail < (spgcnt_t)btopr(vp_len)) {
124         /*
125         * Don't have enough free memory for the
126         * max request, try sizing down vp request.
127         */

```

```

128     deltab = (ssize_t)(off - vp_off);
129     vp_len -= deltab;
130     vp_off += deltab;
131     if (pagesavail < btopr(vp_len)) {
132         /*
133          * Still not enough memory, just settle for
134          * pagesavail which is at least 1.
135          */
136         vp_len = ptob(pagesavail);
137     }
138 }

140 vp_end = vp_off + vp_len;
141 ASSERT(off >= vp_off && off < vp_end);

143 if (isra && segop_kluster(seg, addr, 0))
143 if (isra && SEGOP_KLUSTER(seg, addr, 0))
144     return ((page_t *)NULL); /* segment driver says no */

146 if ((plist = page_create_va(vp, off,
147     PAGESIZE, PG_EXCL | PG_WAIT, seg, addr)) == NULL)
148     return ((page_t *)NULL);

150 if (vp_len <= PAGESIZE || pvn_nofodklust) {
151     *offp = off;
152     *lenp = MIN(vp_len, PAGESIZE);
153 } else {
154     /*
155     * Scan back from front by incrementing "deltab" and
156     * comparing "off" with "vp_off + deltab" to avoid
157     * "signed" versus "unsigned" conversion problems.
158     */
159     for (deltab = PAGESIZE; off >= vp_off + deltab;
160         deltab += PAGESIZE) {
161         /*
162          * Call back to the segment driver to verify that
163          * the klustering/read ahead operation makes sense.
164          */
165         if (segop_kluster(seg, addr, -deltab))
165         if (SEGOP_KLUSTER(seg, addr, -deltab))
166             break; /* page not eligible */
167         if ((pp = page_create_va(vp, off - deltab,
168             PAGESIZE, PG_EXCL, seg, addr - deltab))
169             == NULL)
170             break; /* already have the page */
171         /*
172          * Add page to front of page list.
173          */
174         page_add(&plist, pp);
175     }
176     deltab -= PAGESIZE;

178     /* scan forward from front */
179     for (deltaf = PAGESIZE; off + deltax < vp_end;
180         deltax += PAGESIZE) {
181         /*
182          * Call back to the segment driver to verify that
183          * the klustering/read ahead operation makes sense.
184          */
185         if (segop_kluster(seg, addr, deltax))
185         if (SEGOP_KLUSTER(seg, addr, deltax))
186             break; /* page not file extension */
187         if ((pp = page_create_va(vp, off + deltax,
188             PAGESIZE, PG_EXCL, seg, addr + deltax))
189             == NULL)
190             break; /* already have page */

```

```

192         /*
193          * Add page to end of page list.
194          */
195         page_add(&plist, pp);
196         plist = plist->p_next;
197     }
198     *offp = off = off - deltab;
199     *lenp = deltax + deltax;
200     ASSERT(off >= vp_off);

202     /*
203     * If we ended up getting more than was actually
204     * requested, retract the returned length to only
205     * reflect what was requested. This might happen
206     * if we were allowed to kluster pages across a
207     * span of (say) 5 frags, and frag size is less
208     * than PAGESIZE. We need a whole number of
209     * pages to contain those frags, but the returned
210     * size should only allow the returned range to
211     * extend as far as the end of the frags.
212     */
213     if ((vp_off + vp_len) < (off + *lenp)) {
214         ASSERT(vp_end > off);
215         *lenp = vp_end - off;
216     }
217 }
218 TRACE_3(TR_FAC_VM, TR_PVN_READ_KLUSTER,
219     "pvn_read_kluster:seg %p addr %x isra %x",
220     seg, addr, isra);
221 return (plist);
222 }

```

_____unchanged_portion_omitted_____

new/usr/src/uts/common/vm/vm_seg.c

1

54548 Tue Nov 24 09:34:49 2015

new/usr/src/uts/common/vm/vm_seg.c

patch lower-case-segops

_____ unchanged_portion_omitted_

```
1637 /*
1638  * Unmap a segment and free it from its associated address space.
1639  * This should be called by anybody who's finished with a whole segment's
1640  * mapping. Just calls segop_unmap() on the whole mapping. It is the
1640  * mapping. Just calls SEGOP_UNMAP() on the whole mapping. It is the
1641  * responsibility of the segment driver to unlink the the segment
1642  * from the address space, and to free public and private data structures
1643  * associated with the segment. (This is typically done by a call to
1644  * seg_free()).
1645  */
1646 void
1647 seg_unmap(struct seg *seg)
1648 {
1649 #ifdef DEBUG
1650     int ret;
1651 #endif /* DEBUG */
1652
1653     ASSERT(seg->s_as && AS_WRITE_HELD(seg->s_as, &seg->s_as->a_lock));
1654
1655     /* Shouldn't have called seg_unmap if mapping isn't yet established */
1656     ASSERT(seg->s_data != NULL);
1657
1658     /* Unmap the whole mapping */
1659 #ifdef DEBUG
1660     ret = segop_unmap(seg, seg->s_base, seg->s_size);
1660     ret = SEGOP_UNMAP(seg, seg->s_base, seg->s_size);
1661     ASSERT(ret == 0);
1662 #else
1663     (void) segop_unmap(seg, seg->s_base, seg->s_size);
1663     SEGOP_UNMAP(seg, seg->s_base, seg->s_size);
1664 #endif /* DEBUG */
1665 }
1666
1667 /*
1668  * Free the segment from its associated as. This should only be called
1669  * if a mapping to the segment has not yet been established (e.g., if
1670  * an error occurs in the middle of doing an as_map when the segment
1671  * has already been partially set up) or if it has already been deleted
1672  * (e.g., from a segment driver unmap routine if the unmap applies to the
1673  * entire segment). If the mapping is currently set up then seg_unmap() should
1674  * be called instead.
1675  */
1676 void
1677 seg_free(struct seg *seg)
1678 {
1679     register struct as *as = seg->s_as;
1680     struct seg *tseg = as_removeseg(as, seg);
1681
1682     ASSERT(tseg == seg);
1683
1684     /*
1685      * If the segment private data field is NULL,
1686      * then segment driver is not attached yet.
1687      */
1688     if (seg->s_data != NULL)
1689         segop_free(seg);
1689         SEGOP_FREE(seg);
1690
1691     mutex_destroy(&seg->s_pmtx);
```

new/usr/src/uts/common/vm/vm_seg.c

2

```
1692     ASSERT(seg->s_phead.p_lnext == &seg->s_phead);
1693     ASSERT(seg->s_phead.p_lprev == &seg->s_phead);
1694     kmem_cache_free(seg_cache, seg);
1695 }
1696
1697 _____ unchanged_portion_omitted_
1698
1699 1857 /*
1858  * General not supported function for segop_inherit
1858  * General not supported function for SEGOP_INHERIT
1859  */
1860 /* ARGSUSED */
1861 int
1862 seg_inherit_notsup(struct seg *seg, caddr_t addr, size_t len, uint_t op)
1863 {
1864     return (ENOTSUP);
1865 }
1866
1867 _____ unchanged_portion_omitted_
```