⦃ 1691c07ff4 ⌄   ...

**src** / **sys** / **kern** / vfs_syscalls.c

laffer1 Sync with freebsd                                          ⟳ History

👥 1 contributor

4368 lines (4034 sloc)  │  94 KB                                         ...

```
1    /*-
2     * Copyright (c) 1989, 1993
3     *      The Regents of the University of California.  All rights reserved.
4     * (c) UNIX System Laboratories, Inc.
5     * All or some portions of this file are derived from material licensed
6     * to the University of California by American Telephone and Telegraph
7     * Co. or Unix System Laboratories, Inc. and are reproduced herein with
8     * the permission of UNIX System Laboratories, Inc.
9     *
10    * Redistribution and use in source and binary forms, with or without
11    * modification, are permitted provided that the following conditions
12    * are met:
13    * 1. Redistributions of source code must retain the above copyright
14    *    notice, this list of conditions and the following disclaimer.
15    * 2. Redistributions in binary form must reproduce the above copyright
16    *    notice, this list of conditions and the following disclaimer in the
17    *    documentation and/or other materials provided with the distribution.
18    * 4. Neither the name of the University nor the names of its contributors
19    *    may be used to endorse or promote products derived from this software
20    *    without specific prior written permission.
21    *
22    * THIS SOFTWARE IS PROVIDED BY THE REGENTS AND CONTRIBUTORS ``AS IS'' AND
23    * ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE
24    * IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE
25    * ARE DISCLAIMED.  IN NO EVENT SHALL THE REGENTS OR CONTRIBUTORS BE LIABLE
26    * FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL
27    * DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS
28    * OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION)
29    * HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT
30    * LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY
31    * OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF
32    * SUCH DAMAGE.
33    *
34    *      @(#)vfs_syscalls.c      8.13 (Berkeley) 4/15/94
35    */
36
37   #include <sys/cdefs.h>
38   __FBSDID("$FreeBSD: stable/11/sys/kern/vfs_syscalls.c 338987 2018-09-27 18:54:41Z gordon $");
39
40   #include "opt_capsicum.h"
41   #include "opt_compat.h"
42   #include "opt_ktrace.h"
43
44   #include <sys/param.h>
45   #include <sys/systm.h>
46   #include <sys/bio.h>
47   #include <sys/buf.h>
48   #include <sys/capsicum.h>
49   #include <sys/disk.h>
50   #include <sys/sysent.h>
51   #include <sys/malloc.h>
52   #include <sys/mount.h>
53   #include <sys/mutex.h>
54   #include <sys/sysproto.h>
55   #include <sys/namei.h>
56   #include <sys/filedesc.h>
57   #include <sys/kernel.h>
58   #include <sys/fcntl.h>
59   #include <sys/file.h>
60   #include <sys/filio.h>
61   #include <sys/limits.h>
62   #include <sys/linker.h>
63   #include <sys/rwlock.h>
64   #include <sys/sdt.h>
65   #include <sys/stat.h>
66   #include <sys/sx.h>
67   #include <sys/unistd.h>
68   #include <sys/vnode.h>
69   #include <sys/priv.h>
70   #include <sys/proc.h>
71   #include <sys/dirent.h>
72   #include <sys/jail.h>
73   #include <sys/syscallsubr.h>
74   #include <sys/sysctl.h>
75   #ifdef KTRACE
76   #include <sys/ktrace.h>
77   #endif
78
```

```c
 79    #include <machine/stdarg.h>
 80
 81    #include <security/audit/audit.h>
 82    #include <security/mac/mac_framework.h>
 83
 84    #include <vm/vm.h>
 85    #include <vm/vm_object.h>
 86    #include <vm/vm_page.h>
 87    #include <vm/uma.h>
 88
 89    #include <ufs/ufs/quota.h>
 90
 91    MALLOC_DEFINE(M_FADVISE, "fadvise", "posix_fadvise(2) information");
 92
 93    SDT_PROVIDER_DEFINE(vfs);
 94    SDT_PROBE_DEFINE2(vfs, , stat, mode, "char *", "int");
 95    SDT_PROBE_DEFINE2(vfs, , stat, reg, "char *", "int");
 96
 97    static int kern_chflagsat(struct thread *td, int fd, const char *path,
 98        enum uio_seg pathseg, u_long flags, int atflag);
 99    static int setfflags(struct thread *td, struct vnode *, u_long);
100    static int getutimes(const struct timeval *, enum uio_seg, struct timespec *);
101    static int getutimens(const struct timespec *, enum uio_seg,
102        struct timespec *, int *);
103    static int setutimes(struct thread *td, struct vnode *,
104        const struct timespec *, int, int);
105    static int vn_access(struct vnode *vp, int user_flags, struct ucred *cred,
106        struct thread *td);
107
108    /*
109     * Sync each mounted filesystem.
110     */
111    #ifndef _SYS_SYSPROTO_H_
112    struct sync_args {
113            int     dummy;
114    };
115    #endif
116    /* ARGSUSED */
117    int
118    sys_sync(struct thread *td, struct sync_args *uap)
119    {
120            struct mount *mp, *nmp;
121            int save;
122
123            mtx_lock(&mountlist_mtx);
124            for (mp = TAILQ_FIRST(&mountlist); mp != NULL; mp = nmp) {
125                    if (vfs_busy(mp, MBF_NOWAIT | MBF_MNTLSTLOCK)) {
126                            nmp = TAILQ_NEXT(mp, mnt_list);
127                            continue;
128                    }
129                    if ((mp->mnt_flag & MNT_RDONLY) == 0 &&
130                        vn_start_write(NULL, &mp, V_NOWAIT) == 0) {
131                            save = curthread_pflags_set(TDP_SYNCIO);
132                            vfs_msync(mp, MNT_NOWAIT);
133                            VFS_SYNC(mp, MNT_NOWAIT);
134                            curthread_pflags_restore(save);
135                            vn_finished_write(mp);
136                    }
137                    mtx_lock(&mountlist_mtx);
138                    nmp = TAILQ_NEXT(mp, mnt_list);
139                    vfs_unbusy(mp);
140            }
141            mtx_unlock(&mountlist_mtx);
142            return (0);
143    }
144
145    /*
146     * Change filesystem quotas.
147     */
148    #ifndef _SYS_SYSPROTO_H_
149    struct quotactl_args {
150            char *path;
151            int cmd;
152            int uid;
153            caddr_t arg;
154    };
155    #endif
156    int
157    sys_quotactl(struct thread *td, struct quotactl_args *uap)
158    {
159            struct mount *mp;
160            struct nameidata nd;
161            int error;
162
163            AUDIT_ARG_CMD(uap->cmd);
164            AUDIT_ARG_UID(uap->uid);
165            if (!prison_allow(td->td_ucred, PR_ALLOW_QUOTAS))
166                    return (EPERM);
167            NDINIT(&nd, LOOKUP, FOLLOW | LOCKLEAF | AUDITVNODE1, UIO_USERSPACE,
168                uap->path, td);
169            if ((error = namei(&nd)) != 0)
170                    return (error);
171            NDFREE(&nd, NDF_ONLY_PNBUF);
172            mp = nd.ni_vp->v_mount;
173            vfs_ref(mp);
174            vput(nd.ni_vp);
175            error = vfs_busy(mp, 0);
176            vfs_rel(mp);
```

```
177        if (error != 0)
178                return (error);
179        error = VFS_QUOTACTL(mp, uap->cmd, uap->uid, uap->arg);
180
181        /*
182         * Since quota on operation typically needs to open quota
183         * file, the Q_QUOTAON handler needs to unbusy the mount point
184         * before calling into namei.  Otherwise, unmount might be
185         * started between two vfs_busy() invocations (first is our,
186         * second is from mount point cross-walk code in lookup()),
187         * causing deadlock.
188         *
189         * Require that Q_QUOTAON handles the vfs_busy() reference on
190         * its own, always returning with ubusied mount point.
191         */
192        if ((uap->cmd >> SUBCMDSHIFT) != Q_QUOTAON &&
193            (uap->cmd >> SUBCMDSHIFT) != Q_QUOTAOFF)
194                vfs_unbusy(mp);
195        return (error);
196 }
197
198 /*
199  * Used by statfs conversion routines to scale the block size up if
200  * necessary so that all of the block counts are <= 'max_size'.  Note
201  * that 'max_size' should be a bitmask, i.e. 2^n - 1 for some non-zero
202  * value of 'n'.
203  */
204 void
205 statfs_scale_blocks(struct statfs *sf, long max_size)
206 {
207        uint64_t count;
208        int shift;
209
210        KASSERT(powerof2(max_size + 1), ("%s: invalid max_size", __func__));
211
212        /*
213         * Attempt to scale the block counts to give a more accurate
214         * overview to userland of the ratio of free space to used
215         * space.  To do this, find the largest block count and compute
216         * a divisor that lets it fit into a signed integer <= max_size.
217         */
218        if (sf->f_bavail < 0)
219                count = -sf->f_bavail;
220        else
221                count = sf->f_bavail;
222        count = MAX(sf->f_blocks, MAX(sf->f_bfree, count));
223        if (count <= max_size)
224                return;
225
226        count >>= flsl(max_size);
227        shift = 0;
228        while (count > 0) {
229                shift++;
230                count >>=1;
231        }
232
233        sf->f_bsize <<= shift;
234        sf->f_blocks >>= shift;
235        sf->f_bfree >>= shift;
236        sf->f_bavail >>= shift;
237 }
238
239 static int
240 kern_do_statfs(struct thread *td, struct mount *mp, struct statfs *buf)
241 {
242        struct statfs *sp;
243        int error;
244
245        if (mp == NULL)
246                return (EBADF);
247        error = vfs_busy(mp, 0);
248        vfs_rel(mp);
249        if (error != 0)
250                return (error);
251 #ifdef MAC
252        error = mac_mount_check_stat(td->td_ucred, mp);
253        if (error != 0)
254                goto out;
255 #endif
256        /*
257         * Set these in case the underlying filesystem fails to do so.
258         */
259        sp = &mp->mnt_stat;
260        sp->f_version = STATFS_VERSION;
261        sp->f_namemax = NAME_MAX;
262        sp->f_flags = mp->mnt_flag & MNT_VISFLAGMASK;
263        error = VFS_STATFS(mp, sp);
264        if (error != 0)
265                goto out;
266        *buf = *sp;
267        if (priv_check(td, PRIV_VFS_GENERATION)) {
268                buf->f_fsid.val[0] = buf->f_fsid.val[1] = 0;
269                prison_enforce_statfs(td->td_ucred, mp, buf);
270        }
271 out:
272        vfs_unbusy(mp);
273        return (error);
274 }
```

```c
275
276    /*
277     * Get filesystem statistics.
278     */
279    #ifndef _SYS_SYSPROTO_H_
280    struct statfs_args {
281            char *path;
282            struct statfs *buf;
283    };
284    #endif
285    int
286    sys_statfs(struct thread *td, struct statfs_args *uap)
287    {
288            struct statfs *sfp;
289            int error;
290
291            sfp = malloc(sizeof(struct statfs), M_STATFS, M_WAITOK);
292            error = kern_statfs(td, uap->path, UIO_USERSPACE, sfp);
293            if (error == 0)
294                    error = copyout(sfp, uap->buf, sizeof(struct statfs));
295            free(sfp, M_STATFS);
296            return (error);
297    }
298
299    int
300    kern_statfs(struct thread *td, char *path, enum uio_seg pathseg,
301        struct statfs *buf)
302    {
303            struct mount *mp;
304            struct nameidata nd;
305            int error;
306
307            NDINIT(&nd, LOOKUP, FOLLOW | LOCKSHARED | LOCKLEAF | AUDITVNODE1,
308                pathseg, path, td);
309            error = namei(&nd);
310            if (error != 0)
311                    return (error);
312            mp = nd.ni_vp->v_mount;
313            vfs_ref(mp);
314            NDFREE(&nd, NDF_ONLY_PNBUF);
315            vput(nd.ni_vp);
316            return (kern_do_statfs(td, mp, buf));
317    }
318
319    /*
320     * Get filesystem statistics.
321     */
322    #ifndef _SYS_SYSPROTO_H_
323    struct fstatfs_args {
324            int fd;
325            struct statfs *buf;
326    };
327    #endif
328    int
329    sys_fstatfs(struct thread *td, struct fstatfs_args *uap)
330    {
331            struct statfs *sfp;
332            int error;
333
334            sfp = malloc(sizeof(struct statfs), M_STATFS, M_WAITOK);
335            error = kern_fstatfs(td, uap->fd, sfp);
336            if (error == 0)
337                    error = copyout(sfp, uap->buf, sizeof(struct statfs));
338            free(sfp, M_STATFS);
339            return (error);
340    }
341
342    int
343    kern_fstatfs(struct thread *td, int fd, struct statfs *buf)
344    {
345            struct file *fp;
346            struct mount *mp;
347            struct vnode *vp;
348            cap_rights_t rights;
349            int error;
350
351            AUDIT_ARG_FD(fd);
352            error = getvnode(td, fd, cap_rights_init(&rights, CAP_FSTATFS), &fp);
353            if (error != 0)
354                    return (error);
355            vp = fp->f_vnode;
356            vn_lock(vp, LK_SHARED | LK_RETRY);
357    #ifdef AUDIT
358            AUDIT_ARG_VNODE1(vp);
359    #endif
360            mp = vp->v_mount;
361            if (mp != NULL)
362                    vfs_ref(mp);
363            VOP_UNLOCK(vp, 0);
364            fdrop(fp, td);
365            return (kern_do_statfs(td, mp, buf));
366    }
367
368    /*
369     * Get statistics on all filesystems.
370     */
371    #ifndef _SYS_SYSPROTO_H_
372    struct getfsstat_args {
```

```
373            struct statfs *buf;
374            long bufsize;
375            int mode;
376     };
377     #endif
378     int
379     sys_getfsstat(struct thread *td, struct getfsstat_args *uap)
380     {
381            size_t count;
382            int error;
383
384            if (uap->bufsize < 0 || uap->bufsize > SIZE_MAX)
385                    return (EINVAL);
386            error = kern_getfsstat(td, &uap->buf, uap->bufsize, &count,
387                UIO_USERSPACE, uap->mode);
388            if (error == 0)
389                    td->td_retval[0] = count;
390            return (error);
391     }
392
393     /*
394      * If (bufsize > 0 && bufseg == UIO_SYSSPACE)
395      *      The caller is responsible for freeing memory which will be allocated
396      *      in '*buf'.
397      */
398     int
399     kern_getfsstat(struct thread *td, struct statfs **buf, size_t bufsize,
400         size_t *countp, enum uio_seg bufseg, int mode)
401     {
402            struct mount *mp, *nmp;
403            struct statfs *sfsp, *sp, *sptmp, *tofree;
404            size_t count, maxcount;
405            int error;
406
407            switch (mode) {
408            case MNT_WAIT:
409            case MNT_NOWAIT:
410                    break;
411            default:
412                    return (EINVAL);
413            }
414     restart:
415            maxcount = bufsize / sizeof(struct statfs);
416            if (bufsize == 0) {
417                    sfsp = NULL;
418                    tofree = NULL;
419            } else if (bufseg == UIO_USERSPACE) {
420                    sfsp = *buf;
421                    tofree = NULL;
422            } else /* if (bufseg == UIO_SYSSPACE) */ {
423                    count = 0;
424                    mtx_lock(&mountlist_mtx);
425                    TAILQ_FOREACH(mp, &mountlist, mnt_list) {
426                            count++;
427                    }
428                    mtx_unlock(&mountlist_mtx);
429                    if (maxcount > count)
430                            maxcount = count;
431                    tofree = sfsp = *buf = malloc(maxcount * sizeof(struct statfs),
432                        M_STATFS, M_WAITOK);
433            }
434            count = 0;
435            mtx_lock(&mountlist_mtx);
436            for (mp = TAILQ_FIRST(&mountlist); mp != NULL; mp = nmp) {
437                    if (prison_canseemount(td->td_ucred, mp) != 0) {
438                            nmp = TAILQ_NEXT(mp, mnt_list);
439                            continue;
440                    }
441     #ifdef MAC
442                    if (mac_mount_check_stat(td->td_ucred, mp) != 0) {
443                            nmp = TAILQ_NEXT(mp, mnt_list);
444                            continue;
445                    }
446     #endif
447                    if (mode == MNT_WAIT) {
448                            if (vfs_busy(mp, MBF_MNTLSTLOCK) != 0) {
449                                    /*
450                                     * If vfs_busy() failed, and MBF_NOWAIT
451                                     * wasn't passed, then the mp is gone.
452                                     * Furthermore, because of MBF_MNTLSTLOCK,
453                                     * the mountlist_mtx was dropped.  We have
454                                     * no other choice than to start over.
455                                     */
456                                    mtx_unlock(&mountlist_mtx);
457                                    free(tofree, M_STATFS);
458                                    goto restart;
459                            }
460                    } else {
461                            if (vfs_busy(mp, MBF_NOWAIT | MBF_MNTLSTLOCK) != 0) {
462                                    nmp = TAILQ_NEXT(mp, mnt_list);
463                                    continue;
464                            }
465                    }
466                    if (sfsp != NULL && count < maxcount) {
467                            sp = &mp->mnt_stat;
468                            /*
469                             * Set these in case the underlying filesystem
470                             * fails to do so.
```

```
                           */
                   sp->f_version = STATFS_VERSION;
                   sp->f_namemax = NAME_MAX;
                   sp->f_flags = mp->mnt_flag & MNT_VISFLAGMASK;
                   /*
                    * If MNT_NOWAIT is specified, do not refresh
                    * the fsstat cache.
                    */
                   if (mode != MNT_NOWAIT) {
                           error = VFS_STATFS(mp, sp);
                           if (error != 0) {
                                   mtx_lock(&mountlist_mtx);
                                   nmp = TAILQ_NEXT(mp, mnt_list);
                                   vfs_unbusy(mp);
                                   continue;
                           }
                   }
                   if (priv_check(td, PRIV_VFS_GENERATION)) {
                           sptmp = malloc(sizeof(struct statfs), M_STATFS,
                               M_WAITOK);
                           *sptmp = *sp;
                           sptmp->f_fsid.val[0] = sptmp->f_fsid.val[1] = 0;
                           prison_enforce_statfs(td->td_ucred, mp, sptmp);
                           sp = sptmp;
                   } else
                           sptmp = NULL;
                   if (bufseg == UIO_SYSSPACE) {
                           bcopy(sp, sfsp, sizeof(*sp));
                           free(sptmp, M_STATFS);
                   } else /* if (bufseg == UIO_USERSPACE) */ {
                           error = copyout(sp, sfsp, sizeof(*sp));
                           free(sptmp, M_STATFS);
                           if (error != 0) {
                                   vfs_unbusy(mp);
                                   return (error);
                           }
                   }
                   sfsp++;
           }
           count++;
           mtx_lock(&mountlist_mtx);
           nmp = TAILQ_NEXT(mp, mnt_list);
           vfs_unbusy(mp);
   }
   mtx_unlock(&mountlist_mtx);
   if (sfsp != NULL && count > maxcount)
           *countp = maxcount;
   else
           *countp = count;
   return (0);
}

#ifdef COMPAT_FREEBSD4
/*
 * Get old format filesystem statistics.
 */
static void cvtstatfs(struct statfs *, struct ostatfs *);

#ifndef _SYS_SYSPROTO_H_
struct freebsd4_statfs_args {
        char *path;
        struct ostatfs *buf;
};
#endif
int
freebsd4_statfs(struct thread *td, struct freebsd4_statfs_args *uap)
{
        struct ostatfs osb;
        struct statfs *sfp;
        int error;

        sfp = malloc(sizeof(struct statfs), M_STATFS, M_WAITOK);
        error = kern_statfs(td, uap->path, UIO_USERSPACE, sfp);
        if (error == 0) {
                cvtstatfs(sfp, &osb);
                error = copyout(&osb, uap->buf, sizeof(osb));
        }
        free(sfp, M_STATFS);
        return (error);
}

/*
 * Get filesystem statistics.
 */
#ifndef _SYS_SYSPROTO_H_
struct freebsd4_fstatfs_args {
        int fd;
        struct ostatfs *buf;
};
#endif
int
freebsd4_fstatfs(struct thread *td, struct freebsd4_fstatfs_args *uap)
{
        struct ostatfs osb;
        struct statfs *sfp;
        int error;

        sfp = malloc(sizeof(struct statfs), M_STATFS, M_WAITOK);
```

```
569          error = kern_fstatfs(td, uap->fd, sfp);
570          if (error == 0) {
571                  cvtstatfs(sfp, &osb);
572                  error = copyout(&osb, uap->buf, sizeof(osb));
573          }
574          free(sfp, M_STATFS);
575          return (error);
576  }
577
578  /*
579   * Get statistics on all filesystems.
580   */
581  #ifndef _SYS_SYSPROTO_H_
582  struct freebsd4_getfsstat_args {
583          struct ostatfs *buf;
584          long bufsize;
585          int mode;
586  };
587  #endif
588  int
589  freebsd4_getfsstat(struct thread *td, struct freebsd4_getfsstat_args *uap)
590  {
591          struct statfs *buf, *sp;
592          struct ostatfs osb;
593          size_t count, size;
594          int error;
595
596          if (uap->bufsize < 0)
597                  return (EINVAL);
598          count = uap->bufsize / sizeof(struct ostatfs);
599          if (count > SIZE_MAX / sizeof(struct statfs))
600                  return (EINVAL);
601          size = count * sizeof(struct statfs);
602          error = kern_getfsstat(td, &buf, size, &count, UIO_SYSSPACE,
603              uap->mode);
604          if (buf == NULL)
605                  return (EINVAL);
606          td->td_retval[0] = count;
607          if (size != 0) {
608                  sp = buf;
609                  while (count != 0 && error == 0) {
610                          cvtstatfs(sp, &osb);
611                          error = copyout(&osb, uap->buf, sizeof(osb));
612                          sp++;
613                          uap->buf++;
614                          count--;
615                  }
616                  free(buf, M_STATFS);
617          }
618          return (error);
619  }
620
621  /*
622   * Implement fstatfs() for (NFS) file handles.
623   */
624  #ifndef _SYS_SYSPROTO_H_
625  struct freebsd4_fhstatfs_args {
626          struct fhandle *u_fhp;
627          struct ostatfs *buf;
628  };
629  #endif
630  int
631  freebsd4_fhstatfs(struct thread *td, struct freebsd4_fhstatfs_args *uap)
632  {
633          struct ostatfs osb;
634          struct statfs *sfp;
635          fhandle_t fh;
636          int error;
637
638          error = copyin(uap->u_fhp, &fh, sizeof(fhandle_t));
639          if (error != 0)
640                  return (error);
641          sfp = malloc(sizeof(struct statfs), M_STATFS, M_WAITOK);
642          error = kern_fhstatfs(td, fh, sfp);
643          if (error == 0) {
644                  cvtstatfs(sfp, &osb);
645                  error = copyout(&osb, uap->buf, sizeof(osb));
646          }
647          free(sfp, M_STATFS);
648          return (error);
649  }
650
651  /*
652   * Convert a new format statfs structure to an old format statfs structure.
653   */
654  static void
655  cvtstatfs(struct statfs *nsp, struct ostatfs *osp)
656  {
657
658          statfs_scale_blocks(nsp, LONG_MAX);
659          bzero(osp, sizeof(*osp));
660          osp->f_bsize = nsp->f_bsize;
661          osp->f_iosize = MIN(nsp->f_iosize, LONG_MAX);
662          osp->f_blocks = nsp->f_blocks;
663          osp->f_bfree = nsp->f_bfree;
664          osp->f_bavail = nsp->f_bavail;
665          osp->f_files = MIN(nsp->f_files, LONG_MAX);
666          osp->f_ffree = MIN(nsp->f_ffree, LONG_MAX);
```

```c
667            osp->f_owner = nsp->f_owner;
668            osp->f_type = nsp->f_type;
669            osp->f_flags = nsp->f_flags;
670            osp->f_syncwrites = MIN(nsp->f_syncwrites, LONG_MAX);
671            osp->f_asyncwrites = MIN(nsp->f_asyncwrites, LONG_MAX);
672            osp->f_syncreads = MIN(nsp->f_syncreads, LONG_MAX);
673            osp->f_asyncreads = MIN(nsp->f_asyncreads, LONG_MAX);
674            strlcpy(osp->f_fstypename, nsp->f_fstypename,
675                MIN(MFSNAMELEN, OMFSNAMELEN));
676            strlcpy(osp->f_mntonname, nsp->f_mntonname,
677                MIN(MNAMELEN, OMNAMELEN));
678            strlcpy(osp->f_mntfromname, nsp->f_mntfromname,
679                MIN(MNAMELEN, OMNAMELEN));
680            osp->f_fsid = nsp->f_fsid;
681    }
682    #endif /* COMPAT_FREEBSD4 */
683
684    /*
685     * Change current working directory to a given file descriptor.
686     */
687    #ifndef _SYS_SYSPROTO_H_
688    struct fchdir_args {
689            int     fd;
690    };
691    #endif
692    int
693    sys_fchdir(struct thread *td, struct fchdir_args *uap)
694    {
695            struct vnode *vp, *tdp;
696            struct mount *mp;
697            struct file *fp;
698            cap_rights_t rights;
699            int error;
700
701            AUDIT_ARG_FD(uap->fd);
702            error = getvnode(td, uap->fd, cap_rights_init(&rights, CAP_FCHDIR),
703                &fp);
704            if (error != 0)
705                    return (error);
706            vp = fp->f_vnode;
707            vrefact(vp);
708            fdrop(fp, td);
709            vn_lock(vp, LK_SHARED | LK_RETRY);
710            AUDIT_ARG_VNODE1(vp);
711            error = change_dir(vp, td);
712            while (!error && (mp = vp->v_mountedhere) != NULL) {
713                    if (vfs_busy(mp, 0))
714                            continue;
715                    error = VFS_ROOT(mp, LK_SHARED, &tdp);
716                    vfs_unbusy(mp);
717                    if (error != 0)
718                            break;
719                    vput(vp);
720                    vp = tdp;
721            }
722            if (error != 0) {
723                    vput(vp);
724                    return (error);
725            }
726            VOP_UNLOCK(vp, 0);
727            pwd_chdir(td, vp);
728            return (0);
729    }
730
731    /*
732     * Change current working directory (``.'').
733     */
734    #ifndef _SYS_SYSPROTO_H_
735    struct chdir_args {
736            char    *path;
737    };
738    #endif
739    int
740    sys_chdir(struct thread *td, struct chdir_args *uap)
741    {
742
743            return (kern_chdir(td, uap->path, UIO_USERSPACE));
744    }
745
746    int
747    kern_chdir(struct thread *td, char *path, enum uio_seg pathseg)
748    {
749            struct nameidata nd;
750            int error;
751
752            NDINIT(&nd, LOOKUP, FOLLOW | LOCKSHARED | LOCKLEAF | AUDITVNODE1,
753                pathseg, path, td);
754            if ((error = namei(&nd)) != 0)
755                    return (error);
756            if ((error = change_dir(nd.ni_vp, td)) != 0) {
757                    vput(nd.ni_vp);
758                    NDFREE(&nd, NDF_ONLY_PNBUF);
759                    return (error);
760            }
761            VOP_UNLOCK(nd.ni_vp, 0);
762            NDFREE(&nd, NDF_ONLY_PNBUF);
763            pwd_chdir(td, nd.ni_vp);
764            return (0);
```

```c
765  }
766
767  /*
768   * Change notion of root (``/'') directory.
769   */
770  #ifndef _SYS_SYSPROTO_H_
771  struct chroot_args {
772          char    *path;
773  };
774  #endif
775  int
776  sys_chroot(struct thread *td, struct chroot_args *uap)
777  {
778          struct nameidata nd;
779          int error;
780
781          error = priv_check(td, PRIV_VFS_CHROOT);
782          if (error != 0)
783                  return (error);
784          NDINIT(&nd, LOOKUP, FOLLOW | LOCKSHARED | LOCKLEAF | AUDITVNODE1,
785              UIO_USERSPACE, uap->path, td);
786          error = namei(&nd);
787          if (error != 0)
788                  goto error;
789          error = change_dir(nd.ni_vp, td);
790          if (error != 0)
791                  goto e_vunlock;
792  #ifdef MAC
793          error = mac_vnode_check_chroot(td->td_ucred, nd.ni_vp);
794          if (error != 0)
795                  goto e_vunlock;
796  #endif
797          VOP_UNLOCK(nd.ni_vp, 0);
798          error = pwd_chroot(td, nd.ni_vp);
799          vrele(nd.ni_vp);
800          NDFREE(&nd, NDF_ONLY_PNBUF);
801          return (error);
802  e_vunlock:
803          vput(nd.ni_vp);
804  error:
805          NDFREE(&nd, NDF_ONLY_PNBUF);
806          return (error);
807  }
808
809  /*
810   * Common routine for chroot and chdir.  Callers must provide a locked vnode
811   * instance.
812   */
813  int
814  change_dir(struct vnode *vp, struct thread *td)
815  {
816  #ifdef MAC
817          int error;
818  #endif
819
820          ASSERT_VOP_LOCKED(vp, "change_dir(): vp not locked");
821          if (vp->v_type != VDIR)
822                  return (ENOTDIR);
823  #ifdef MAC
824          error = mac_vnode_check_chdir(td->td_ucred, vp);
825          if (error != 0)
826                  return (error);
827  #endif
828          return (VOP_ACCESS(vp, VEXEC, td->td_ucred, td));
829  }
830
831  static __inline void
832  flags_to_rights(int flags, cap_rights_t *rightsp)
833  {
834
835          if (flags & O_EXEC) {
836                  cap_rights_set(rightsp, CAP_FEXECVE);
837          } else {
838                  switch ((flags & O_ACCMODE)) {
839                  case O_RDONLY:
840                          cap_rights_set(rightsp, CAP_READ);
841                          break;
842                  case O_RDWR:
843                          cap_rights_set(rightsp, CAP_READ);
844                          /* FALLTHROUGH */
845                  case O_WRONLY:
846                          cap_rights_set(rightsp, CAP_WRITE);
847                          if (!(flags & (O_APPEND | O_TRUNC)))
848                                  cap_rights_set(rightsp, CAP_SEEK);
849                          break;
850                  }
851          }
852
853          if (flags & O_CREAT)
854                  cap_rights_set(rightsp, CAP_CREATE);
855
856          if (flags & O_TRUNC)
857                  cap_rights_set(rightsp, CAP_FTRUNCATE);
858
859          if (flags & (O_SYNC | O_FSYNC))
860                  cap_rights_set(rightsp, CAP_FSYNC);
861
862          if (flags & (O_EXLOCK | O_SHLOCK))
```

```c
863                        cap_rights_set(rightsp, CAP_FLOCK);
864    }
865
866    /*
867     * Check permissions, allocate an open file structure, and call the device
868     * open routine if any.
869     */
870    #ifndef _SYS_SYSPROTO_H_
871    struct open_args {
872            char    *path;
873            int     flags;
874            int     mode;
875    };
876    #endif
877    int
878    sys_open(struct thread *td, struct open_args *uap)
879    {
880
881            return (kern_openat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
882                uap->flags, uap->mode));
883    }
884
885    #ifndef _SYS_SYSPROTO_H_
886    struct openat_args {
887            int     fd;
888            char    *path;
889            int     flag;
890            int     mode;
891    };
892    #endif
893    int
894    sys_openat(struct thread *td, struct openat_args *uap)
895    {
896
897            AUDIT_ARG_FD(uap->fd);
898            return (kern_openat(td, uap->fd, uap->path, UIO_USERSPACE, uap->flag,
899                uap->mode));
900    }
901
902    int
903    kern_openat(struct thread *td, int fd, char *path, enum uio_seg pathseg,
904        int flags, int mode)
905    {
906            struct proc *p = td->td_proc;
907            struct filedesc *fdp = p->p_fd;
908            struct file *fp;
909            struct vnode *vp;
910            struct nameidata nd;
911            cap_rights_t rights;
912            int cmode, error, indx;
913
914            indx = -1;
915
916            AUDIT_ARG_FFLAGS(flags);
917            AUDIT_ARG_MODE(mode);
918            cap_rights_init(&rights, CAP_LOOKUP);
919            flags_to_rights(flags, &rights);
920            /*
921             * Only one of the O_EXEC, O_RDONLY, O_WRONLY and O_RDWR flags
922             * may be specified.
923             */
924            if (flags & O_EXEC) {
925                    if (flags & O_ACCMODE)
926                            return (EINVAL);
927            } else if ((flags & O_ACCMODE) == O_ACCMODE) {
928                    return (EINVAL);
929            } else {
930                    flags = FFLAGS(flags);
931            }
932
933            /*
934             * Allocate a file structure. The descriptor to reference it
935             * is allocated and set by finstall() below.
936             */
937            error = falloc_noinstall(td, &fp);
938            if (error != 0)
939                    return (error);
940            /*
941             * An extra reference on `fp' has been held for us by
942             * falloc_noinstall().
943             */
944            /* Set the flags early so the finit in devfs can pick them up. */
945            fp->f_flag = flags & FMASK;
946            cmode = ((mode & ~fdp->fd_cmask) & ALLPERMS) & ~S_ISTXT;
947            NDINIT_ATRIGHTS(&nd, LOOKUP, FOLLOW | AUDITVNODE1, pathseg, path, fd,
948                &rights, td);
949            td->td_dupfd = -1;              /* XXX check for fdopen */
950            error = vn_open(&nd, &flags, cmode, fp);
951            if (error != 0) {
952                    /*
953                     * If the vn_open replaced the method vector, something
954                     * wonderous happened deep below and we just pass it up
955                     * pretending we know what we do.
956                     */
957                    if (error == ENXIO && fp->f_ops != &badfileops)
958                            goto success;
959
960                    /*
```

```
961                    * Handle special fdopen() case. bleh.
962                    *
963                    * Don't do this for relative (capability) lookups; we don't
964                    * understand exactly what would happen, and we don't think
965                    * that it ever should.
966                    */
967                   if ((nd.ni_lcf & NI_LCF_STRICTRELATIVE) == 0 &&
968                       (error == ENODEV || error == ENXIO) &&
969                       td->td_dupfd >= 0) {
970                           error = dupfdopen(td, fdp, td->td_dupfd, flags, error,
971                               &indx);
972                           if (error == 0)
973                                   goto success;
974                   }
975
976                   goto bad;
977           }
978           td->td_dupfd = 0;
979           NDFREE(&nd, NDF_ONLY_PNBUF);
980           vp = nd.ni_vp;
981
982           /*
983            * Store the vnode, for any f_type. Typically, the vnode use
984            * count is decremented by direct call to vn_closefile() for
985            * files that switched type in the cdevsw fdopen() method.
986            */
987           fp->f_vnode = vp;
988           /*
989            * If the file wasn't claimed by devfs bind it to the normal
990            * vnode operations here.
991            */
992           if (fp->f_ops == &badfileops) {
993                   KASSERT(vp->v_type != VFIFO, ("Unexpected fifo."));
994                   fp->f_seqcount = 1;
995                   finit(fp, (flags & FMASK) | (fp->f_flag & FHASLOCK),
996                       DTYPE_VNODE, vp, &vnops);
997           }
998
999           VOP_UNLOCK(vp, 0);
1000          if (flags & O_TRUNC) {
1001                  error = fo_truncate(fp, 0, td->td_ucred, td);
1002                  if (error != 0)
1003                          goto bad;
1004          }
1005   success:
1006          /*
1007           * If we haven't already installed the FD (for dupfdopen), do so now.
1008           */
1009          if (indx == -1) {
1010                  struct filecaps *fcaps;
1011
1012   #ifdef CAPABILITIES
1013                  if ((nd.ni_lcf & NI_LCF_STRICTRELATIVE) != 0)
1014                          fcaps = &nd.ni_filecaps;
1015                  else
1016   #endif
1017                          fcaps = NULL;
1018                  error = finstall(td, fp, &indx, flags, fcaps);
1019                  /* On success finstall() consumes fcaps. */
1020                  if (error != 0) {
1021                          filecaps_free(&nd.ni_filecaps);
1022                          goto bad;
1023                  }
1024          } else {
1025                  filecaps_free(&nd.ni_filecaps);
1026          }
1027
1028          /*
1029           * Release our private reference, leaving the one associated with
1030           * the descriptor table intact.
1031           */
1032          fdrop(fp, td);
1033          td->td_retval[0] = indx;
1034          return (0);
1035   bad:
1036          KASSERT(indx == -1, ("indx=%d, should be -1", indx));
1037          fdrop(fp, td);
1038          return (error);
1039   }
1040
1041   #ifdef COMPAT_43
1042   /*
1043    * Create a file.
1044    */
1045   #ifndef _SYS_SYSPROTO_H_
1046   struct ocreat_args {
1047           char    *path;
1048           int     mode;
1049   };
1050   #endif
1051   int
1052   ocreat(struct thread *td, struct ocreat_args *uap)
1053   {
1054
1055          return (kern_openat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
1056              O_WRONLY | O_CREAT | O_TRUNC, uap->mode));
1057   }
1058   #endif /* COMPAT_43 */
```

```c
/*
 * Create a special file.
 */
#ifndef _SYS_SYSPROTO_H_
struct mknod_args {
        char    *path;
        int     mode;
        int     dev;
};
#endif
int
sys_mknod(struct thread *td, struct mknod_args *uap)
{

        return (kern_mknodat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
            uap->mode, uap->dev));
}


#ifndef _SYS_SYSPROTO_H_
struct mknodat_args {
        int     fd;
        char    *path;
        mode_t  mode;
        dev_t   dev;
};
#endif
int
sys_mknodat(struct thread *td, struct mknodat_args *uap)
{

        return (kern_mknodat(td, uap->fd, uap->path, UIO_USERSPACE, uap->mode,
            uap->dev));
}

int
kern_mknodat(struct thread *td, int fd, char *path, enum uio_seg pathseg,
    int mode, int dev)
{
        struct vnode *vp;
        struct mount *mp;
        struct vattr vattr;
        struct nameidata nd;
        cap_rights_t rights;
        int error, whiteout = 0;

        AUDIT_ARG_MODE(mode);
        AUDIT_ARG_DEV(dev);
        switch (mode & S_IFMT) {
        case S_IFCHR:
        case S_IFBLK:
                error = priv_check(td, PRIV_VFS_MKNOD_DEV);
                if (error == 0 && dev == VNOVAL)
                        error = EINVAL;
                break;
        case S_IFWHT:
                error = priv_check(td, PRIV_VFS_MKNOD_WHT);
                break;
        case S_IFIFO:
                if (dev == 0)
                        return (kern_mkfifoat(td, fd, path, pathseg, mode));
                /* FALLTHROUGH */
        default:
                error = EINVAL;
                break;
        }
        if (error != 0)
                return (error);
restart:
        bwillwrite();
        NDINIT_ATRIGHTS(&nd, CREATE, LOCKPARENT | SAVENAME | AUDITVNODE1 |
            NOCACHE, pathseg, path, fd, cap_rights_init(&rights, CAP_MKNODAT),
            td);
        if ((error = namei(&nd)) != 0)
                return (error);
        vp = nd.ni_vp;
        if (vp != NULL) {
                NDFREE(&nd, NDF_ONLY_PNBUF);
                if (vp == nd.ni_dvp)
                        vrele(nd.ni_dvp);
                else
                        vput(nd.ni_dvp);
                vrele(vp);
                return (EEXIST);
        } else {
                VATTR_NULL(&vattr);
                vattr.va_mode = (mode & ALLPERMS) &
                    ~td->td_proc->p_fd->fd_cmask;
                vattr.va_rdev = dev;
                whiteout = 0;

                switch (mode & S_IFMT) {
                case S_IFCHR:
                        vattr.va_type = VCHR;
                        break;
                case S_IFBLK:
                        vattr.va_type = VBLK;
                        break;
```

```c
                case S_IFWHT:
                        whiteout = 1;
                        break;
                default:
                        panic("kern_mknod: invalid mode");
                }
        }
        if (vn_start_write(nd.ni_dvp, &mp, V_NOWAIT) != 0) {
                NDFREE(&nd, NDF_ONLY_PNBUF);
                vput(nd.ni_dvp);
                if ((error = vn_start_write(NULL, &mp, V_XSLEEP | PCATCH)) != 0)
                        return (error);
                goto restart;
        }
#ifdef MAC
        if (error == 0 && !whiteout)
                error = mac_vnode_check_create(td->td_ucred, nd.ni_dvp,
                    &nd.ni_cnd, &vattr);
#endif
        if (error == 0) {
                if (whiteout)
                        error = VOP_WHITEOUT(nd.ni_dvp, &nd.ni_cnd, CREATE);
                else {
                        error = VOP_MKNOD(nd.ni_dvp, &nd.ni_vp,
                                                &nd.ni_cnd, &vattr);
                        if (error == 0)
                                vput(nd.ni_vp);
                }
        }
        NDFREE(&nd, NDF_ONLY_PNBUF);
        vput(nd.ni_dvp);
        vn_finished_write(mp);
        return (error);
}

/*
 * Create a named pipe.
 */
#ifndef _SYS_SYSPROTO_H_
struct mkfifo_args {
        char    *path;
        int     mode;
};
#endif
int
sys_mkfifo(struct thread *td, struct mkfifo_args *uap)
{

        return (kern_mkfifoat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
            uap->mode));
}

#ifndef _SYS_SYSPROTO_H_
struct mkfifoat_args {
        int     fd;
        char    *path;
        mode_t  mode;
};
#endif
int
sys_mkfifoat(struct thread *td, struct mkfifoat_args *uap)
{

        return (kern_mkfifoat(td, uap->fd, uap->path, UIO_USERSPACE,
            uap->mode));
}

int
kern_mkfifoat(struct thread *td, int fd, char *path, enum uio_seg pathseg,
    int mode)
{
        struct mount *mp;
        struct vattr vattr;
        struct nameidata nd;
        cap_rights_t rights;
        int error;

        AUDIT_ARG_MODE(mode);
restart:
        bwillwrite();
        NDINIT_ATRIGHTS(&nd, CREATE, LOCKPARENT | SAVENAME | AUDITVNODE1 |
            NOCACHE, pathseg, path, fd, cap_rights_init(&rights, CAP_MKFIFOAT),
            td);
        if ((error = namei(&nd)) != 0)
                return (error);
        if (nd.ni_vp != NULL) {
                NDFREE(&nd, NDF_ONLY_PNBUF);
                if (nd.ni_vp == nd.ni_dvp)
                        vrele(nd.ni_dvp);
                else
                        vput(nd.ni_dvp);
                vrele(nd.ni_vp);
                return (EEXIST);
        }
        if (vn_start_write(nd.ni_dvp, &mp, V_NOWAIT) != 0) {
                NDFREE(&nd, NDF_ONLY_PNBUF);
                vput(nd.ni_dvp);
                if ((error = vn_start_write(NULL, &mp, V_XSLEEP | PCATCH)) != 0)
```

```c
                            return (error);
                    goto restart;
            }
            VATTR_NULL(&vattr);
            vattr.va_type = VFIFO;
            vattr.va_mode = (mode & ALLPERMS) & ~td->td_proc->p_fd->fd_cmask;
#ifdef MAC
            error = mac_vnode_check_create(td->td_ucred, nd.ni_dvp, &nd.ni_cnd,
                &vattr);
            if (error != 0)
                    goto out;
#endif
            error = VOP_MKNOD(nd.ni_dvp, &nd.ni_vp, &nd.ni_cnd, &vattr);
            if (error == 0)
                    vput(nd.ni_vp);
#ifdef MAC
out:
#endif
            vput(nd.ni_dvp);
            vn_finished_write(mp);
            NDFREE(&nd, NDF_ONLY_PNBUF);
            return (error);
}

/*
 * Make a hard file link.
 */
#ifndef _SYS_SYSPROTO_H_
struct link_args {
        char    *path;
        char    *link;
};
#endif
int
sys_link(struct thread *td, struct link_args *uap)
{

        return (kern_linkat(td, AT_FDCWD, AT_FDCWD, uap->path, uap->link,
            UIO_USERSPACE, FOLLOW));
}

#ifndef _SYS_SYSPROTO_H_
struct linkat_args {
        int     fd1;
        char    *path1;
        int     fd2;
        char    *path2;
        int     flag;
};
#endif
int
sys_linkat(struct thread *td, struct linkat_args *uap)
{
        int flag;

        flag = uap->flag;
        if (flag & ~AT_SYMLINK_FOLLOW)
                return (EINVAL);

        return (kern_linkat(td, uap->fd1, uap->fd2, uap->path1, uap->path2,
            UIO_USERSPACE, (flag & AT_SYMLINK_FOLLOW) ? FOLLOW : NOFOLLOW));
}

int hardlink_check_uid = 0;
SYSCTL_INT(_security_bsd, OID_AUTO, hardlink_check_uid, CTLFLAG_RW,
    &hardlink_check_uid, 0,
    "Unprivileged processes cannot create hard links to files owned by other "
    "users");
static int hardlink_check_gid = 0;
SYSCTL_INT(_security_bsd, OID_AUTO, hardlink_check_gid, CTLFLAG_RW,
    &hardlink_check_gid, 0,
    "Unprivileged processes cannot create hard links to files owned by other "
    "groups");

static int
can_hardlink(struct vnode *vp, struct ucred *cred)
{
        struct vattr va;
        int error;

        if (!hardlink_check_uid && !hardlink_check_gid)
                return (0);

        error = VOP_GETATTR(vp, &va, cred);
        if (error != 0)
                return (error);

        if (hardlink_check_uid && cred->cr_uid != va.va_uid) {
                error = priv_check_cred(cred, PRIV_VFS_LINK, 0);
                if (error != 0)
                        return (error);
        }

        if (hardlink_check_gid && !groupmember(va.va_gid, cred)) {
                error = priv_check_cred(cred, PRIV_VFS_LINK, 0);
                if (error != 0)
                        return (error);
        }
```

```c
1353
1354            return (0);
1355    }
1356
1357    int
1358    kern_linkat(struct thread *td, int fd1, int fd2, char *path1, char *path2,
1359        enum uio_seg segflg, int follow)
1360    {
1361            struct vnode *vp;
1362            struct mount *mp;
1363            struct nameidata nd;
1364            cap_rights_t rights;
1365            int error;
1366
1367    again:
1368            bwillwrite();
1369            NDINIT_ATRIGHTS(&nd, LOOKUP, follow | AUDITVNODE1, segflg, path1, fd1,
1370                cap_rights_init(&rights, CAP_LINKAT_SOURCE), td);
1371
1372            if ((error = namei(&nd)) != 0)
1373                    return (error);
1374            NDFREE(&nd, NDF_ONLY_PNBUF);
1375            vp = nd.ni_vp;
1376            if (vp->v_type == VDIR) {
1377                    vrele(vp);
1378                    return (EPERM);         /* POSIX */
1379            }
1380            NDINIT_ATRIGHTS(&nd, CREATE,
1381                LOCKPARENT | SAVENAME | AUDITVNODE2 | NOCACHE, segflg, path2, fd2,
1382                cap_rights_init(&rights, CAP_LINKAT_TARGET), td);
1383            if ((error = namei(&nd)) == 0) {
1384                    if (nd.ni_vp != NULL) {
1385                            NDFREE(&nd, NDF_ONLY_PNBUF);
1386                            if (nd.ni_dvp == nd.ni_vp)
1387                                    vrele(nd.ni_dvp);
1388                            else
1389                                    vput(nd.ni_dvp);
1390                            vrele(nd.ni_vp);
1391                            vrele(vp);
1392                            return (EEXIST);
1393                    } else if (nd.ni_dvp->v_mount != vp->v_mount) {
1394                            /*
1395                             * Cross-device link.  No need to recheck
1396                             * vp->v_type, since it cannot change, except
1397                             * to VBAD.
1398                             */
1399                            NDFREE(&nd, NDF_ONLY_PNBUF);
1400                            vput(nd.ni_dvp);
1401                            vrele(vp);
1402                            return (EXDEV);
1403                    } else if ((error = vn_lock(vp, LK_EXCLUSIVE)) == 0) {
1404                            error = can_hardlink(vp, td->td_ucred);
1405    #ifdef MAC
1406                            if (error == 0)
1407                                    error = mac_vnode_check_link(td->td_ucred,
1408                                        nd.ni_dvp, vp, &nd.ni_cnd);
1409    #endif
1410                            if (error != 0) {
1411                                    vput(vp);
1412                                    vput(nd.ni_dvp);
1413                                    NDFREE(&nd, NDF_ONLY_PNBUF);
1414                                    return (error);
1415                            }
1416                            error = vn_start_write(vp, &mp, V_NOWAIT);
1417                            if (error != 0) {
1418                                    vput(vp);
1419                                    vput(nd.ni_dvp);
1420                                    NDFREE(&nd, NDF_ONLY_PNBUF);
1421                                    error = vn_start_write(NULL, &mp,
1422                                        V_XSLEEP | PCATCH);
1423                                    if (error != 0)
1424                                            return (error);
1425                                    goto again;
1426                            }
1427                            error = VOP_LINK(nd.ni_dvp, vp, &nd.ni_cnd);
1428                            VOP_UNLOCK(vp, 0);
1429                            vput(nd.ni_dvp);
1430                            vn_finished_write(mp);
1431                            NDFREE(&nd, NDF_ONLY_PNBUF);
1432                    } else {
1433                            vput(nd.ni_dvp);
1434                            NDFREE(&nd, NDF_ONLY_PNBUF);
1435                            vrele(vp);
1436                            goto again;
1437                    }
1438            }
1439            vrele(vp);
1440            return (error);
1441    }
1442
1443    /*
1444     * Make a symbolic link.
1445     */
1446    #ifndef _SYS_SYSPROTO_H_
1447    struct symlink_args {
1448            char    *path;
1449            char    *link;
1450    };
```

```
1451    #endif
1452    int
1453    sys_symlink(struct thread *td, struct symlink_args *uap)
1454    {
1455
1456            return (kern_symlinkat(td, uap->path, AT_FDCWD, uap->link,
1457                UIO_USERSPACE));
1458    }
1459
1460    #ifndef _SYS_SYSPROTO_H_
1461    struct symlinkat_args {
1462            char    *path;
1463            int     fd;
1464            char    *path2;
1465    };
1466    #endif
1467    int
1468    sys_symlinkat(struct thread *td, struct symlinkat_args *uap)
1469    {
1470
1471            return (kern_symlinkat(td, uap->path1, uap->fd, uap->path2,
1472                UIO_USERSPACE));
1473    }
1474
1475    int
1476    kern_symlinkat(struct thread *td, char *path1, int fd, char *path2,
1477        enum uio_seg segflg)
1478    {
1479            struct mount *mp;
1480            struct vattr vattr;
1481            char *syspath;
1482            struct nameidata nd;
1483            int error;
1484            cap_rights_t rights;
1485
1486            if (segflg == UIO_SYSSPACE) {
1487                    syspath = path1;
1488            } else {
1489                    syspath = uma_zalloc(namei_zone, M_WAITOK);
1490                    if ((error = copyinstr(path1, syspath, MAXPATHLEN, NULL)) != 0)
1491                            goto out;
1492            }
1493            AUDIT_ARG_TEXT(syspath);
1494    restart:
1495            bwillwrite();
1496            NDINIT_ATRIGHTS(&nd, CREATE, LOCKPARENT | SAVENAME | AUDITVNODE1 |
1497                NOCACHE, segflg, path2, fd, cap_rights_init(&rights, CAP_SYMLINKAT),
1498                td);
1499            if ((error = namei(&nd)) != 0)
1500                    goto out;
1501            if (nd.ni_vp) {
1502                    NDFREE(&nd, NDF_ONLY_PNBUF);
1503                    if (nd.ni_vp == nd.ni_dvp)
1504                            vrele(nd.ni_dvp);
1505                    else
1506                            vput(nd.ni_dvp);
1507                    vrele(nd.ni_vp);
1508                    error = EEXIST;
1509                    goto out;
1510            }
1511            if (vn_start_write(nd.ni_dvp, &mp, V_NOWAIT) != 0) {
1512                    NDFREE(&nd, NDF_ONLY_PNBUF);
1513                    vput(nd.ni_dvp);
1514                    if ((error = vn_start_write(NULL, &mp, V_XSLEEP | PCATCH)) != 0)
1515                            goto out;
1516                    goto restart;
1517            }
1518            VATTR_NULL(&vattr);
1519            vattr.va_mode = ACCESSPERMS &~ td->td_proc->p_fd->fd_cmask;
1520    #ifdef MAC
1521            vattr.va_type = VLNK;
1522            error = mac_vnode_check_create(td->td_ucred, nd.ni_dvp, &nd.ni_cnd,
1523                &vattr);
1524            if (error != 0)
1525                    goto out2;
1526    #endif
1527            error = VOP_SYMLINK(nd.ni_dvp, &nd.ni_vp, &nd.ni_cnd, &vattr, syspath);
1528            if (error == 0)
1529                    vput(nd.ni_vp);
1530    #ifdef MAC
1531    out2:
1532    #endif
1533            NDFREE(&nd, NDF_ONLY_PNBUF);
1534            vput(nd.ni_dvp);
1535            vn_finished_write(mp);
1536    out:
1537            if (segflg != UIO_SYSSPACE)
1538                    uma_zfree(namei_zone, syspath);
1539            return (error);
1540    }
1541
1542    /*
1543     * Delete a whiteout from the filesystem.
1544     */
1545    #ifndef _SYS_SYSPROTO_H_
1546    struct undelete_args {
1547            char *path;
1548    };
```

```c
#endif
int
sys_undelete(struct thread *td, struct undelete_args *uap)
{
        struct mount *mp;
        struct nameidata nd;
        int error;

restart:
        bwillwrite();
        NDINIT(&nd, DELETE, LOCKPARENT | DOWHITEOUT | AUDITVNODE1,
            UIO_USERSPACE, uap->path, td);
        error = namei(&nd);
        if (error != 0)
                return (error);

        if (nd.ni_vp != NULLVP || !(nd.ni_cnd.cn_flags & ISWHITEOUT)) {
                NDFREE(&nd, NDF_ONLY_PNBUF);
                if (nd.ni_vp == nd.ni_dvp)
                        vrele(nd.ni_dvp);
                else
                        vput(nd.ni_dvp);
                if (nd.ni_vp)
                        vrele(nd.ni_vp);
                return (EEXIST);
        }
        if (vn_start_write(nd.ni_dvp, &mp, V_NOWAIT) != 0) {
                NDFREE(&nd, NDF_ONLY_PNBUF);
                vput(nd.ni_dvp);
                if ((error = vn_start_write(NULL, &mp, V_XSLEEP | PCATCH)) != 0)
                        return (error);
                goto restart;
        }
        error = VOP_WHITEOUT(nd.ni_dvp, &nd.ni_cnd, DELETE);
        NDFREE(&nd, NDF_ONLY_PNBUF);
        vput(nd.ni_dvp);
        vn_finished_write(mp);
        return (error);
}

/*
 * Delete a name from the filesystem.
 */
#ifndef _SYS_SYSPROTO_H_
struct unlink_args {
        char    *path;
};
#endif
int
sys_unlink(struct thread *td, struct unlink_args *uap)
{

        return (kern_unlinkat(td, AT_FDCWD, uap->path, UIO_USERSPACE, 0));
}

#ifndef _SYS_SYSPROTO_H_
struct unlinkat_args {
        int     fd;
        char    *path;
        int     flag;
};
#endif
int
sys_unlinkat(struct thread *td, struct unlinkat_args *uap)
{
        int flag = uap->flag;
        int fd = uap->fd;
        char *path = uap->path;

        if (flag & ~AT_REMOVEDIR)
                return (EINVAL);

        if (flag & AT_REMOVEDIR)
                return (kern_rmdirat(td, fd, path, UIO_USERSPACE));
        else
                return (kern_unlinkat(td, fd, path, UIO_USERSPACE, 0));
}

int
kern_unlinkat(struct thread *td, int fd, char *path, enum uio_seg pathseg,
    ino_t oldinum)
{
        struct mount *mp;
        struct vnode *vp;
        struct nameidata nd;
        struct stat sb;
        cap_rights_t rights;
        int error;

restart:
        bwillwrite();
        NDINIT_ATRIGHTS(&nd, DELETE, LOCKPARENT | LOCKLEAF | AUDITVNODE1,
            pathseg, path, fd, cap_rights_init(&rights, CAP_UNLINKAT), td);
        if ((error = namei(&nd)) != 0)
                return (error == EINVAL ? EPERM : error);
        vp = nd.ni_vp;
        if (vp->v_type == VDIR && oldinum == 0) {
                error = EPERM;          /* POSIX */
```

```c
       } else if (oldinum != 0 &&
               ((error = vn_stat(vp, &sb, td->td_ucred, NOCRED, td)) == 0) &&
               sb.st_ino != oldinum) {
                   error = EIDRM;  /* Identifier removed */
       } else {
           /*
            * The root of a mounted filesystem cannot be deleted.
            *
            * XXX: can this only be a VDIR case?
            */
           if (vp->v_vflag & VV_ROOT)
                   error = EBUSY;
       }
       if (error == 0) {
           if (vn_start_write(nd.ni_dvp, &mp, V_NOWAIT) != 0) {
                   NDFREE(&nd, NDF_ONLY_PNBUF);
                   vput(nd.ni_dvp);
                   if (vp == nd.ni_dvp)
                           vrele(vp);
                   else
                           vput(vp);
                   if ((error = vn_start_write(NULL, &mp,
                       V_XSLEEP | PCATCH)) != 0)
                           return (error);
                   goto restart;
           }
#ifdef MAC
           error = mac_vnode_check_unlink(td->td_ucred, nd.ni_dvp, vp,
               &nd.ni_cnd);
           if (error != 0)
                   goto out;
#endif
           vfs_notify_upper(vp, VFS_NOTIFY_UPPER_UNLINK);
           error = VOP_REMOVE(nd.ni_dvp, vp, &nd.ni_cnd);
#ifdef MAC
out:
#endif
           vn_finished_write(mp);
       }
       NDFREE(&nd, NDF_ONLY_PNBUF);
       vput(nd.ni_dvp);
       if (vp == nd.ni_dvp)
               vrele(vp);
       else
               vput(vp);
       return (error);
}

/*
 * Reposition read/write file offset.
 */
#ifndef _SYS_SYSPROTO_H_
struct lseek_args {
       int     fd;
       int     pad;
       off_t   offset;
       int     whence;
};
#endif
int
sys_lseek(struct thread *td, struct lseek_args *uap)
{

       return (kern_lseek(td, uap->fd, uap->offset, uap->whence));
}

int
kern_lseek(struct thread *td, int fd, off_t offset, int whence)
{
       struct file *fp;
       cap_rights_t rights;
       int error;

       AUDIT_ARG_FD(fd);
       error = fget(td, fd, cap_rights_init(&rights, CAP_SEEK), &fp);
       if (error != 0)
               return (error);
       error = (fp->f_ops->fo_flags & DFLAG_SEEKABLE) != 0 ?
           fo_seek(fp, offset, whence, td) : ESPIPE;
       fdrop(fp, td);
       return (error);
}

#if defined(COMPAT_43)
/*
 * Reposition read/write file offset.
 */
#ifndef _SYS_SYSPROTO_H_
struct olseek_args {
       int     fd;
       long    offset;
       int     whence;
};
#endif
int
olseek(struct thread *td, struct olseek_args *uap)
{
```

```
1745            return (kern_lseek(td, uap->fd, uap->offset, uap->whence));
1746    }
1747    #endif /* COMPAT_43 */

1749    #if defined(COMPAT_FREEBSD6)
1750    /* Version with the 'pad' argument */
1751    int
1752    freebsd6_lseek(struct thread *td, struct freebsd6_lseek_args *uap)
1753    {

1755            return (kern_lseek(td, uap->fd, uap->offset, uap->whence));
1756    }
1757    #endif

1759    /*
1760     * Check access permissions using passed credentials.
1761     */
1762    static int
1763    vn_access(struct vnode *vp, int user_flags, struct ucred *cred,
1764        struct thread *td)
1765    {
1766            accmode_t accmode;
1767            int error;

1769            /* Flags == 0 means only check for existence. */
1770            if (user_flags == 0)
1771                    return (0);

1773            accmode = 0;
1774            if (user_flags & R_OK)
1775                    accmode |= VREAD;
1776            if (user_flags & W_OK)
1777                    accmode |= VWRITE;
1778            if (user_flags & X_OK)
1779                    accmode |= VEXEC;
1780    #ifdef MAC
1781            error = mac_vnode_check_access(cred, vp, accmode);
1782            if (error != 0)
1783                    return (error);
1784    #endif
1785            if ((accmode & VWRITE) == 0 || (error = vn_writechk(vp)) == 0)
1786                    error = VOP_ACCESS(vp, accmode, cred, td);
1787            return (error);
1788    }

1790    /*
1791     * Check access permissions using "real" credentials.
1792     */
1793    #ifndef _SYS_SYSPROTO_H_
1794    struct access_args {
1795            char    *path;
1796            int     amode;
1797    };
1798    #endif
1799    int
1800    sys_access(struct thread *td, struct access_args *uap)
1801    {

1803            return (kern_accessat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
1804                0, uap->amode));
1805    }

1807    #ifndef _SYS_SYSPROTO_H_
1808    struct faccessat_args {
1809            int     dirfd;
1810            char    *path;
1811            int     amode;
1812            int     flag;
1813    }
1814    #endif
1815    int
1816    sys_faccessat(struct thread *td, struct faccessat_args *uap)
1817    {

1819            return (kern_accessat(td, uap->fd, uap->path, UIO_USERSPACE, uap->flag,
1820                uap->amode));
1821    }

1823    int
1824    kern_accessat(struct thread *td, int fd, char *path, enum uio_seg pathseg,
1825        int flag, int amode)
1826    {
1827            struct ucred *cred, *usecred;
1828            struct vnode *vp;
1829            struct nameidata nd;
1830            cap_rights_t rights;
1831            int error;

1833            if (flag & ~AT_EACCESS)
1834                    return (EINVAL);
1835            if (amode != F_OK && (amode & ~(R_OK | W_OK | X_OK)) != 0)
1836                    return (EINVAL);

1838            /*
1839             * Create and modify a temporary credential instead of one that
1840             * is potentially shared (if we need one).
1841             */
1842            cred = td->td_ucred;
```

```c
1843          if ((flag & AT_EACCESS) == 0 &&
1844              ((cred->cr_uid != cred->cr_ruid ||
1845              cred->cr_rgid != cred->cr_groups[0]))) {
1846                  usecred = crdup(cred);
1847                  usecred->cr_uid = cred->cr_ruid;
1848                  usecred->cr_groups[0] = cred->cr_rgid;
1849                  td->td_ucred = usecred;
1850          } else
1851                  usecred = cred;
1852          AUDIT_ARG_VALUE(amode);
1853          NDINIT_ATRIGHTS(&nd, LOOKUP, FOLLOW | LOCKSHARED | LOCKLEAF |
1854              AUDITVNODE1, pathseg, path, fd, cap_rights_init(&rights, CAP_FSTAT),
1855              td);
1856          if ((error = namei(&nd)) != 0)
1857                  goto out;
1858          vp = nd.ni_vp;
1859
1860          error = vn_access(vp, amode, usecred, td);
1861          NDFREE(&nd, NDF_ONLY_PNBUF);
1862          vput(vp);
1863 out:
1864          if (usecred != cred) {
1865                  td->td_ucred = cred;
1866                  crfree(usecred);
1867          }
1868          return (error);
1869 }
1870
1871 /*
1872  * Check access permissions using "effective" credentials.
1873  */
1874 #ifndef _SYS_SYSPROTO_H_
1875 struct eaccess_args {
1876          char    *path;
1877          int     amode;
1878 };
1879 #endif
1880 int
1881 sys_eaccess(struct thread *td, struct eaccess_args *uap)
1882 {
1883
1884          return (kern_accessat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
1885              AT_EACCESS, uap->amode));
1886 }
1887
1888 #if defined(COMPAT_43)
1889 /*
1890  * Get file status; this version follows links.
1891  */
1892 #ifndef _SYS_SYSPROTO_H_
1893 struct ostat_args {
1894          char    *path;
1895          struct ostat *ub;
1896 };
1897 #endif
1898 int
1899 ostat(struct thread *td, struct ostat_args *uap)
1900 {
1901          struct stat sb;
1902          struct ostat osb;
1903          int error;
1904
1905          error = kern_statat(td, 0, AT_FDCWD, uap->path, UIO_USERSPACE,
1906              &sb, NULL);
1907          if (error != 0)
1908                  return (error);
1909          cvtstat(&sb, &osb);
1910          return (copyout(&osb, uap->ub, sizeof (osb)));
1911 }
1912
1913 /*
1914  * Get file status; this version does not follow links.
1915  */
1916 #ifndef _SYS_SYSPROTO_H_
1917 struct olstat_args {
1918          char    *path;
1919          struct ostat *ub;
1920 };
1921 #endif
1922 int
1923 olstat(struct thread *td, struct olstat_args *uap)
1924 {
1925          struct stat sb;
1926          struct ostat osb;
1927          int error;
1928
1929          error = kern_statat(td, AT_SYMLINK_NOFOLLOW, AT_FDCWD, uap->path,
1930              UIO_USERSPACE, &sb, NULL);
1931          if (error != 0)
1932                  return (error);
1933          cvtstat(&sb, &osb);
1934          return (copyout(&osb, uap->ub, sizeof (osb)));
1935 }
1936
1937 /*
1938  * Convert from an old to a new stat structure.
1939  */
1940 void
```

```c
cvtstat(struct stat *st, struct ostat *ost)
{

	bzero(ost, sizeof(*ost));
	ost->st_dev = st->st_dev;
	ost->st_ino = st->st_ino;
	ost->st_mode = st->st_mode;
	ost->st_nlink = st->st_nlink;
	ost->st_uid = st->st_uid;
	ost->st_gid = st->st_gid;
	ost->st_rdev = st->st_rdev;
	if (st->st_size < (quad_t)1 << 32)
		ost->st_size = st->st_size;
	else
		ost->st_size = -2;
	ost->st_atim = st->st_atim;
	ost->st_mtim = st->st_mtim;
	ost->st_ctim = st->st_ctim;
	ost->st_blksize = st->st_blksize;
	ost->st_blocks = st->st_blocks;
	ost->st_flags = st->st_flags;
	ost->st_gen = st->st_gen;
}
#endif /* COMPAT_43 */

/*
 * Get file status; this version follows links.
 */
#ifndef _SYS_SYSPROTO_H_
struct stat_args {
	char	*path;
	struct stat *ub;
};
#endif
int
sys_stat(struct thread *td, struct stat_args *uap)
{
	struct stat sb;
	int error;

	error = kern_statat(td, 0, AT_FDCWD, uap->path, UIO_USERSPACE,
	    &sb, NULL);
	if (error == 0)
		error = copyout(&sb, uap->ub, sizeof (sb));
	return (error);
}

#ifndef _SYS_SYSPROTO_H_
struct fstatat_args {
	int	fd;
	char	*path;
	struct stat	*buf;
	int	flag;
}
#endif
int
sys_fstatat(struct thread *td, struct fstatat_args *uap)
{
	struct stat sb;
	int error;

	error = kern_statat(td, uap->flag, uap->fd, uap->path,
	    UIO_USERSPACE, &sb, NULL);
	if (error == 0)
		error = copyout(&sb, uap->buf, sizeof (sb));
	return (error);
}

int
kern_statat(struct thread *td, int flag, int fd, char *path,
    enum uio_seg pathseg, struct stat *sbp,
    void (*hook)(struct vnode *vp, struct stat *sbp))
{
	struct nameidata nd;
	struct stat sb;
	cap_rights_t rights;
	int error;

	if (flag & ~AT_SYMLINK_NOFOLLOW)
		return (EINVAL);

	NDINIT_ATRIGHTS(&nd, LOOKUP, ((flag & AT_SYMLINK_NOFOLLOW) ? NOFOLLOW :
	    FOLLOW) | LOCKSHARED | LOCKLEAF | AUDITVNODE1, pathseg, path, fd,
	    cap_rights_init(&rights, CAP_FSTAT), td);

	if ((error = namei(&nd)) != 0)
		return (error);
	error = vn_stat(nd.ni_vp, &sb, td->td_ucred, NOCRED, td);
	if (error == 0) {
		SDT_PROBE2(vfs, , stat, mode, path, sb.st_mode);
		if (S_ISREG(sb.st_mode))
			SDT_PROBE2(vfs, , stat, reg, path, pathseg);
		if (__predict_false(hook != NULL))
			hook(nd.ni_vp, &sb);
	}
	NDFREE(&nd, NDF_ONLY_PNBUF);
	vput(nd.ni_vp);
	if (error != 0)
```

```
2039                        return (error);
2040                *sbp = sb;
2041    #ifdef KTRACE
2042                if (KTRPOINT(td, KTR_STRUCT))
2043                        ktrstat(&sb);
2044    #endif
2045                return (0);
2046    }
2047
2048    /*
2049     * Get file status; this version does not follow links.
2050     */
2051    #ifndef _SYS_SYSPROTO_H_
2052    struct lstat_args {
2053                char    *path;
2054                struct stat *ub;
2055    };
2056    #endif
2057    int
2058    sys_lstat(struct thread *td, struct lstat_args *uap)
2059    {
2060                struct stat sb;
2061                int error;
2062
2063                error = kern_statat(td, AT_SYMLINK_NOFOLLOW, AT_FDCWD, uap->path,
2064                    UIO_USERSPACE, &sb, NULL);
2065                if (error == 0)
2066                        error = copyout(&sb, uap->ub, sizeof (sb));
2067                return (error);
2068    }
2069
2070    /*
2071     * Implementation of the NetBSD [l]stat() functions.
2072     */
2073    void
2074    cvtnstat( struct stat *sb, struct nstat *nsb)
2075    {
2076
2077                bzero(nsb, sizeof *nsb);
2078                nsb->st_dev = sb->st_dev;
2079                nsb->st_ino = sb->st_ino;
2080                nsb->st_mode = sb->st_mode;
2081                nsb->st_nlink = sb->st_nlink;
2082                nsb->st_uid = sb->st_uid;
2083                nsb->st_gid = sb->st_gid;
2084                nsb->st_rdev = sb->st_rdev;
2085                nsb->st_atim = sb->st_atim;
2086                nsb->st_mtim = sb->st_mtim;
2087                nsb->st_ctim = sb->st_ctim;
2088                nsb->st_size = sb->st_size;
2089                nsb->st_blocks = sb->st_blocks;
2090                nsb->st_blksize = sb->st_blksize;
2091                nsb->st_flags = sb->st_flags;
2092                nsb->st_gen = sb->st_gen;
2093                nsb->st_birthtim = sb->st_birthtim;
2094    }
2095
2096    #ifndef _SYS_SYSPROTO_H_
2097    struct nstat_args {
2098                char    *path;
2099                struct nstat *ub;
2100    };
2101    #endif
2102    int
2103    sys_nstat(struct thread *td, struct nstat_args *uap)
2104    {
2105                struct stat sb;
2106                struct nstat nsb;
2107                int error;
2108
2109                error = kern_statat(td, 0, AT_FDCWD, uap->path, UIO_USERSPACE,
2110                    &sb, NULL);
2111                if (error != 0)
2112                        return (error);
2113                cvtnstat(&sb, &nsb);
2114                return (copyout(&nsb, uap->ub, sizeof (nsb)));
2115    }
2116
2117    /*
2118     * NetBSD lstat.  Get file status; this version does not follow links.
2119     */
2120    #ifndef _SYS_SYSPROTO_H_
2121    struct lstat_args {
2122                char    *path;
2123                struct stat *ub;
2124    };
2125    #endif
2126    int
2127    sys_nlstat(struct thread *td, struct nlstat_args *uap)
2128    {
2129                struct stat sb;
2130                struct nstat nsb;
2131                int error;
2132
2133                error = kern_statat(td, AT_SYMLINK_NOFOLLOW, AT_FDCWD, uap->path,
2134                    UIO_USERSPACE, &sb, NULL);
2135                if (error != 0)
2136                        return (error);
```

```c
2137            cvtnstat(&sb, &nsb);
2138            return (copyout(&nsb, uap->ub, sizeof (nsb)));
2139    }
2140
2141    /*
2142     * Get configurable pathname variables.
2143     */
2144    #ifndef _SYS_SYSPROTO_H_
2145    struct pathconf_args {
2146            char    *path;
2147            int     name;
2148    };
2149    #endif
2150    int
2151    sys_pathconf(struct thread *td, struct pathconf_args *uap)
2152    {
2153
2154            return (kern_pathconf(td, uap->path, UIO_USERSPACE, uap->name, FOLLOW));
2155    }
2156
2157    #ifndef _SYS_SYSPROTO_H_
2158    struct lpathconf_args {
2159            char    *path;
2160            int     name;
2161    };
2162    #endif
2163    int
2164    sys_lpathconf(struct thread *td, struct lpathconf_args *uap)
2165    {
2166
2167            return (kern_pathconf(td, uap->path, UIO_USERSPACE, uap->name,
2168                NOFOLLOW));
2169    }
2170
2171    int
2172    kern_pathconf(struct thread *td, char *path, enum uio_seg pathseg, int name,
2173        u_long flags)
2174    {
2175            struct nameidata nd;
2176            int error;
2177
2178            NDINIT(&nd, LOOKUP, LOCKSHARED | LOCKLEAF | AUDITVNODE1 | flags,
2179                pathseg, path, td);
2180            if ((error = namei(&nd)) != 0)
2181                    return (error);
2182            NDFREE(&nd, NDF_ONLY_PNBUF);
2183
2184            error = VOP_PATHCONF(nd.ni_vp, name, td->td_retval);
2185            vput(nd.ni_vp);
2186            return (error);
2187    }
2188
2189    /*
2190     * Return target name of a symbolic link.
2191     */
2192    #ifndef _SYS_SYSPROTO_H_
2193    struct readlink_args {
2194            char    *path;
2195            char    *buf;
2196            size_t  count;
2197    };
2198    #endif
2199    int
2200    sys_readlink(struct thread *td, struct readlink_args *uap)
2201    {
2202
2203            return (kern_readlinkat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
2204                uap->buf, UIO_USERSPACE, uap->count));
2205    }
2206    #ifndef _SYS_SYSPROTO_H_
2207    struct readlinkat_args {
2208            int     fd;
2209            char    *path;
2210            char    *buf;
2211            size_t  bufsize;
2212    };
2213    #endif
2214    int
2215    sys_readlinkat(struct thread *td, struct readlinkat_args *uap)
2216    {
2217
2218            return (kern_readlinkat(td, uap->fd, uap->path, UIO_USERSPACE,
2219                uap->buf, UIO_USERSPACE, uap->bufsize));
2220    }
2221
2222    int
2223    kern_readlinkat(struct thread *td, int fd, char *path, enum uio_seg pathseg,
2224        char *buf, enum uio_seg bufseg, size_t count)
2225    {
2226            struct vnode *vp;
2227            struct iovec aiov;
2228            struct uio auio;
2229            struct nameidata nd;
2230            int error;
2231
2232            if (count > IOSIZE_MAX)
2233                    return (EINVAL);
2234
```

```
2235        NDINIT_AT(&nd, LOOKUP, NOFOLLOW | LOCKSHARED | LOCKLEAF | AUDITVNODE1,
2236            pathseg, path, fd, td);
2237
2238        if ((error = namei(&nd)) != 0)
2239            return (error);
2240        NDFREE(&nd, NDF_ONLY_PNBUF);
2241        vp = nd.ni_vp;
2242    #ifdef MAC
2243        error = mac_vnode_check_readlink(td->td_ucred, vp);
2244        if (error != 0) {
2245            vput(vp);
2246            return (error);
2247        }
2248    #endif
2249        if (vp->v_type != VLNK && (vp->v_vflag & VV_READLINK) == 0)
2250            error = EINVAL;
2251        else {
2252            aiov.iov_base = buf;
2253            aiov.iov_len = count;
2254            auio.uio_iov = &aiov;
2255            auio.uio_iovcnt = 1;
2256            auio.uio_offset = 0;
2257            auio.uio_rw = UIO_READ;
2258            auio.uio_segflg = bufseg;
2259            auio.uio_td = td;
2260            auio.uio_resid = count;
2261            error = VOP_READLINK(vp, &auio, td->td_ucred);
2262            td->td_retval[0] = count - auio.uio_resid;
2263        }
2264        vput(vp);
2265        return (error);
2266    }
2267
2268    /*
2269     * Common implementation code for chflags() and fchflags().
2270     */
2271    static int
2272    setfflags(struct thread *td, struct vnode *vp, u_long flags)
2273    {
2274        struct mount *mp;
2275        struct vattr vattr;
2276        int error;
2277
2278        /* We can't support the value matching VNOVAL. */
2279        if (flags == VNOVAL)
2280            return (EOPNOTSUPP);
2281
2282        /*
2283         * Prevent non-root users from setting flags on devices.  When
2284         * a device is reused, users can retain ownership of the device
2285         * if they are allowed to set flags and programs assume that
2286         * chown can't fail when done as root.
2287         */
2288        if (vp->v_type == VCHR || vp->v_type == VBLK) {
2289            error = priv_check(td, PRIV_VFS_CHFLAGS_DEV);
2290            if (error != 0)
2291                return (error);
2292        }
2293
2294        if ((error = vn_start_write(vp, &mp, V_WAIT | PCATCH)) != 0)
2295            return (error);
2296        VATTR_NULL(&vattr);
2297        vattr.va_flags = flags;
2298        vn_lock(vp, LK_EXCLUSIVE | LK_RETRY);
2299    #ifdef MAC
2300        error = mac_vnode_check_setflags(td->td_ucred, vp, vattr.va_flags);
2301        if (error == 0)
2302    #endif
2303            error = VOP_SETATTR(vp, &vattr, td->td_ucred);
2304        VOP_UNLOCK(vp, 0);
2305        vn_finished_write(mp);
2306        return (error);
2307    }
2308
2309    /*
2310     * Change flags of a file given a path name.
2311     */
2312    #ifndef _SYS_SYSPROTO_H_
2313    struct chflags_args {
2314        const char *path;
2315        u_long  flags;
2316    };
2317    #endif
2318    int
2319    sys_chflags(struct thread *td, struct chflags_args *uap)
2320    {
2321
2322        return (kern_chflagsat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
2323            uap->flags, 0));
2324    }
2325
2326    #ifndef _SYS_SYSPROTO_H_
2327    struct chflagsat_args {
2328        int     fd;
2329        const char *path;
2330        u_long  flags;
2331        int     atflag;
2332    }
```

```c
#endif
int
sys_chflagsat(struct thread *td, struct chflagsat_args *uap)
{
        int fd = uap->fd;
        const char *path = uap->path;
        u_long flags = uap->flags;
        int atflag = uap->atflag;

        if (atflag & ~AT_SYMLINK_NOFOLLOW)
                return (EINVAL);

        return (kern_chflagsat(td, fd, path, UIO_USERSPACE, flags, atflag));
}

/*
 * Same as chflags() but doesn't follow symlinks.
 */
#ifndef _SYS_SYSPROTO_H_
struct lchflags_args {
        const char *path;
        u_long flags;
};
#endif
int
sys_lchflags(struct thread *td, struct lchflags_args *uap)
{

        return (kern_chflagsat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
            uap->flags, AT_SYMLINK_NOFOLLOW));
}

static int
kern_chflagsat(struct thread *td, int fd, const char *path,
    enum uio_seg pathseg, u_long flags, int atflag)
{
        struct nameidata nd;
        cap_rights_t rights;
        int error, follow;

        AUDIT_ARG_FFLAGS(flags);
        follow = (atflag & AT_SYMLINK_NOFOLLOW) ? NOFOLLOW : FOLLOW;
        NDINIT_ATRIGHTS(&nd, LOOKUP, follow | AUDITVNODE1, pathseg, path, fd,
            cap_rights_init(&rights, CAP_FCHFLAGS), td);
        if ((error = namei(&nd)) != 0)
                return (error);
        NDFREE(&nd, NDF_ONLY_PNBUF);
        error = setfflags(td, nd.ni_vp, flags);
        vrele(nd.ni_vp);
        return (error);
}

/*
 * Change flags of a file given a file descriptor.
 */
#ifndef _SYS_SYSPROTO_H_
struct fchflags_args {
        int     fd;
        u_long  flags;
};
#endif
int
sys_fchflags(struct thread *td, struct fchflags_args *uap)
{
        struct file *fp;
        cap_rights_t rights;
        int error;

        AUDIT_ARG_FD(uap->fd);
        AUDIT_ARG_FFLAGS(uap->flags);
        error = getvnode(td, uap->fd, cap_rights_init(&rights, CAP_FCHFLAGS),
            &fp);
        if (error != 0)
                return (error);
#ifdef AUDIT
        vn_lock(fp->f_vnode, LK_SHARED | LK_RETRY);
        AUDIT_ARG_VNODE1(fp->f_vnode);
        VOP_UNLOCK(fp->f_vnode, 0);
#endif
        error = setfflags(td, fp->f_vnode, uap->flags);
        fdrop(fp, td);
        return (error);
}

/*
 * Common implementation code for chmod(), lchmod() and fchmod().
 */
int
setfmode(struct thread *td, struct ucred *cred, struct vnode *vp, int mode)
{
        struct mount *mp;
        struct vattr vattr;
        int error;

        if ((error = vn_start_write(vp, &mp, V_WAIT | PCATCH)) != 0)
                return (error);
        vn_lock(vp, LK_EXCLUSIVE | LK_RETRY);
        VATTR_NULL(&vattr);
```

```c
2431            vattr.va_mode = mode & ALLPERMS;
2432 #ifdef MAC
2433            error = mac_vnode_check_setmode(cred, vp, vattr.va_mode);
2434            if (error == 0)
2435 #endif
2436                    error = VOP_SETATTR(vp, &vattr, cred);
2437            VOP_UNLOCK(vp, 0);
2438            vn_finished_write(mp);
2439            return (error);
2440 }
2441
2442 /*
2443  * Change mode of a file given path name.
2444  */
2445 #ifndef _SYS_SYSPROTO_H_
2446 struct chmod_args {
2447            char    *path;
2448            int     mode;
2449 };
2450 #endif
2451 int
2452 sys_chmod(struct thread *td, struct chmod_args *uap)
2453 {
2454
2455            return (kern_fchmodat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
2456                uap->mode, 0));
2457 }
2458
2459 #ifndef _SYS_SYSPROTO_H_
2460 struct fchmodat_args {
2461            int     dirfd;
2462            char    *path;
2463            mode_t  mode;
2464            int     flag;
2465 }
2466 #endif
2467 int
2468 sys_fchmodat(struct thread *td, struct fchmodat_args *uap)
2469 {
2470            int flag = uap->flag;
2471            int fd = uap->fd;
2472            char *path = uap->path;
2473            mode_t mode = uap->mode;
2474
2475            if (flag & ~AT_SYMLINK_NOFOLLOW)
2476                    return (EINVAL);
2477
2478            return (kern_fchmodat(td, fd, path, UIO_USERSPACE, mode, flag));
2479 }
2480
2481 /*
2482  * Change mode of a file given path name (don't follow links.)
2483  */
2484 #ifndef _SYS_SYSPROTO_H_
2485 struct lchmod_args {
2486            char    *path;
2487            int     mode;
2488 };
2489 #endif
2490 int
2491 sys_lchmod(struct thread *td, struct lchmod_args *uap)
2492 {
2493
2494            return (kern_fchmodat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
2495                uap->mode, AT_SYMLINK_NOFOLLOW));
2496 }
2497
2498 int
2499 kern_fchmodat(struct thread *td, int fd, char *path, enum uio_seg pathseg,
2500     mode_t mode, int flag)
2501 {
2502            struct nameidata nd;
2503            cap_rights_t rights;
2504            int error, follow;
2505
2506            AUDIT_ARG_MODE(mode);
2507            follow = (flag & AT_SYMLINK_NOFOLLOW) ? NOFOLLOW : FOLLOW;
2508            NDINIT_ATRIGHTS(&nd, LOOKUP, follow | AUDITVNODE1, pathseg, path, fd,
2509                cap_rights_init(&rights, CAP_FCHMOD), td);
2510            if ((error = namei(&nd)) != 0)
2511                    return (error);
2512            NDFREE(&nd, NDF_ONLY_PNBUF);
2513            error = setfmode(td, td->td_ucred, nd.ni_vp, mode);
2514            vrele(nd.ni_vp);
2515            return (error);
2516 }
2517
2518 /*
2519  * Change mode of a file given a file descriptor.
2520  */
2521 #ifndef _SYS_SYSPROTO_H_
2522 struct fchmod_args {
2523            int     fd;
2524            int     mode;
2525 };
2526 #endif
2527 int
2528 sys_fchmod(struct thread *td, struct fchmod_args *uap)
```

```
2529   {
2530           struct file *fp;
2531           cap_rights_t rights;
2532           int error;
2533
2534           AUDIT_ARG_FD(uap->fd);
2535           AUDIT_ARG_MODE(uap->mode);
2536
2537           error = fget(td, uap->fd, cap_rights_init(&rights, CAP_FCHMOD), &fp);
2538           if (error != 0)
2539                   return (error);
2540           error = fo_chmod(fp, uap->mode, td->td_ucred, td);
2541           fdrop(fp, td);
2542           return (error);
2543   }
2544
2545   /*
2546    * Common implementation for chown(), lchown(), and fchown()
2547    */
2548   int
2549   setfown(struct thread *td, struct ucred *cred, struct vnode *vp, uid_t uid,
2550       gid_t gid)
2551   {
2552           struct mount *mp;
2553           struct vattr vattr;
2554           int error;
2555
2556           if ((error = vn_start_write(vp, &mp, V_WAIT | PCATCH)) != 0)
2557                   return (error);
2558           vn_lock(vp, LK_EXCLUSIVE | LK_RETRY);
2559           VATTR_NULL(&vattr);
2560           vattr.va_uid = uid;
2561           vattr.va_gid = gid;
2562   #ifdef MAC
2563           error = mac_vnode_check_setowner(cred, vp, vattr.va_uid,
2564               vattr.va_gid);
2565           if (error == 0)
2566   #endif
2567                   error = VOP_SETATTR(vp, &vattr, cred);
2568           VOP_UNLOCK(vp, 0);
2569           vn_finished_write(mp);
2570           return (error);
2571   }
2572
2573   /*
2574    * Set ownership given a path name.
2575    */
2576   #ifndef _SYS_SYSPROTO_H_
2577   struct chown_args {
2578           char    *path;
2579           int     uid;
2580           int     gid;
2581   };
2582   #endif
2583   int
2584   sys_chown(struct thread *td, struct chown_args *uap)
2585   {
2586
2587           return (kern_fchownat(td, AT_FDCWD, uap->path, UIO_USERSPACE, uap->uid,
2588               uap->gid, 0));
2589   }
2590
2591   #ifndef _SYS_SYSPROTO_H_
2592   struct fchownat_args {
2593           int fd;
2594           const char * path;
2595           uid_t uid;
2596           gid_t gid;
2597           int flag;
2598   };
2599   #endif
2600   int
2601   sys_fchownat(struct thread *td, struct fchownat_args *uap)
2602   {
2603           int flag;
2604
2605           flag = uap->flag;
2606           if (flag & ~AT_SYMLINK_NOFOLLOW)
2607                   return (EINVAL);
2608
2609           return (kern_fchownat(td, uap->fd, uap->path, UIO_USERSPACE, uap->uid,
2610               uap->gid, uap->flag));
2611   }
2612
2613   int
2614   kern_fchownat(struct thread *td, int fd, char *path, enum uio_seg pathseg,
2615       int uid, int gid, int flag)
2616   {
2617           struct nameidata nd;
2618           cap_rights_t rights;
2619           int error, follow;
2620
2621           AUDIT_ARG_OWNER(uid, gid);
2622           follow = (flag & AT_SYMLINK_NOFOLLOW) ? NOFOLLOW : FOLLOW;
2623           NDINIT_ATRIGHTS(&nd, LOOKUP, follow | AUDITVNODE1, pathseg, path, fd,
2624               cap_rights_init(&rights, CAP_FCHOWN), td);
2625
2626           if ((error = namei(&nd)) != 0)
```

```
2627                    return (error);
2628            NDFREE(&nd, NDF_ONLY_PNBUF);
2629            error = setfown(td, td->td_ucred, nd.ni_vp, uid, gid);
2630            vrele(nd.ni_vp);
2631            return (error);
2632    }
2633
2634    /*
2635     * Set ownership given a path name, do not cross symlinks.
2636     */
2637    #ifndef _SYS_SYSPROTO_H_
2638    struct lchown_args {
2639            char    *path;
2640            int     uid;
2641            int     gid;
2642    };
2643    #endif
2644    int
2645    sys_lchown(struct thread *td, struct lchown_args *uap)
2646    {
2647
2648            return (kern_fchownat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
2649                uap->uid, uap->gid, AT_SYMLINK_NOFOLLOW));
2650    }
2651
2652    /*
2653     * Set ownership given a file descriptor.
2654     */
2655    #ifndef _SYS_SYSPROTO_H_
2656    struct fchown_args {
2657            int     fd;
2658            int     uid;
2659            int     gid;
2660    };
2661    #endif
2662    int
2663    sys_fchown(struct thread *td, struct fchown_args *uap)
2664    {
2665            struct file *fp;
2666            cap_rights_t rights;
2667            int error;
2668
2669            AUDIT_ARG_FD(uap->fd);
2670            AUDIT_ARG_OWNER(uap->uid, uap->gid);
2671            error = fget(td, uap->fd, cap_rights_init(&rights, CAP_FCHOWN), &fp);
2672            if (error != 0)
2673                    return (error);
2674            error = fo_chown(fp, uap->uid, uap->gid, td->td_ucred, td);
2675            fdrop(fp, td);
2676            return (error);
2677    }
2678
2679    /*
2680     * Common implementation code for utimes(), lutimes(), and futimes().
2681     */
2682    static int
2683    getutimes(const struct timeval *usrtvp, enum uio_seg tvpseg,
2684        struct timespec *tsp)
2685    {
2686            struct timeval tv[2];
2687            const struct timeval *tvp;
2688            int error;
2689
2690            if (usrtvp == NULL) {
2691                    vfs_timestamp(&tsp[0]);
2692                    tsp[1] = tsp[0];
2693            } else {
2694                    if (tvpseg == UIO_SYSSPACE) {
2695                            tvp = usrtvp;
2696                    } else {
2697                            if ((error = copyin(usrtvp, tv, sizeof(tv))) != 0)
2698                                    return (error);
2699                            tvp = tv;
2700                    }
2701
2702                    if (tvp[0].tv_usec < 0 || tvp[0].tv_usec >= 1000000 ||
2703                        tvp[1].tv_usec < 0 || tvp[1].tv_usec >= 1000000)
2704                            return (EINVAL);
2705                    TIMEVAL_TO_TIMESPEC(&tvp[0], &tsp[0]);
2706                    TIMEVAL_TO_TIMESPEC(&tvp[1], &tsp[1]);
2707            }
2708            return (0);
2709    }
2710
2711    /*
2712     * Common implementation code for futimens(), utimensat().
2713     */
2714    #define UTIMENS_NULL    0x1
2715    #define UTIMENS_EXIT    0x2
2716    static int
2717    getutimens(const struct timespec *usrtsp, enum uio_seg tspseg,
2718        struct timespec *tsp, int *retflags)
2719    {
2720            struct timespec tsnow;
2721            int error;
2722
2723            vfs_timestamp(&tsnow);
2724            *retflags = 0;
```

```
2725        if (usrtsp == NULL) {
2726                tsp[0] = tsnow;
2727                tsp[1] = tsnow;
2728                *retflags |= UTIMENS_NULL;
2729                return (0);
2730        }
2731        if (tspseg == UIO_SYSSPACE) {
2732                tsp[0] = usrtsp[0];
2733                tsp[1] = usrtsp[1];
2734        } else if ((error = copyin(usrtsp, tsp, sizeof(*tsp) * 2)) != 0)
2735                return (error);
2736        if (tsp[0].tv_nsec == UTIME_OMIT && tsp[1].tv_nsec == UTIME_OMIT)
2737                *retflags |= UTIMENS_EXIT;
2738        if (tsp[0].tv_nsec == UTIME_NOW && tsp[1].tv_nsec == UTIME_NOW)
2739                *retflags |= UTIMENS_NULL;
2740        if (tsp[0].tv_nsec == UTIME_OMIT)
2741                tsp[0].tv_sec = VNOVAL;
2742        else if (tsp[0].tv_nsec == UTIME_NOW)
2743                tsp[0] = tsnow;
2744        else if (tsp[0].tv_nsec < 0 || tsp[0].tv_nsec >= 1000000000L)
2745                return (EINVAL);
2746        if (tsp[1].tv_nsec == UTIME_OMIT)
2747                tsp[1].tv_sec = VNOVAL;
2748        else if (tsp[1].tv_nsec == UTIME_NOW)
2749                tsp[1] = tsnow;
2750        else if (tsp[1].tv_nsec < 0 || tsp[1].tv_nsec >= 1000000000L)
2751                return (EINVAL);
2752
2753        return (0);
2754 }
2755
2756 /*
2757  * Common implementation code for utimes(), lutimes(), futimes(), futimens(),
2758  * and utimensat().
2759  */
2760 static int
2761 setutimes(struct thread *td, struct vnode *vp, const struct timespec *ts,
2762     int numtimes, int nullflag)
2763 {
2764        struct mount *mp;
2765        struct vattr vattr;
2766        int error, setbirthtime;
2767
2768        if ((error = vn_start_write(vp, &mp, V_WAIT | PCATCH)) != 0)
2769                return (error);
2770        vn_lock(vp, LK_EXCLUSIVE | LK_RETRY);
2771        setbirthtime = 0;
2772        if (numtimes < 3 && !VOP_GETATTR(vp, &vattr, td->td_ucred) &&
2773            timespeccmp(&ts[1], &vattr.va_birthtime, < ))
2774                setbirthtime = 1;
2775        VATTR_NULL(&vattr);
2776        vattr.va_atime = ts[0];
2777        vattr.va_mtime = ts[1];
2778        if (setbirthtime)
2779                vattr.va_birthtime = ts[1];
2780        if (numtimes > 2)
2781                vattr.va_birthtime = ts[2];
2782        if (nullflag)
2783                vattr.va_vaflags |= VA_UTIMES_NULL;
2784 #ifdef MAC
2785        error = mac_vnode_check_setutimes(td->td_ucred, vp, vattr.va_atime,
2786            vattr.va_mtime);
2787 #endif
2788        if (error == 0)
2789                error = VOP_SETATTR(vp, &vattr, td->td_ucred);
2790        VOP_UNLOCK(vp, 0);
2791        vn_finished_write(mp);
2792        return (error);
2793 }
2794
2795 /*
2796  * Set the access and modification times of a file.
2797  */
2798 #ifndef _SYS_SYSPROTO_H_
2799 struct utimes_args {
2800        char    *path;
2801        struct  timeval *tptr;
2802 };
2803 #endif
2804 int
2805 sys_utimes(struct thread *td, struct utimes_args *uap)
2806 {
2807
2808        return (kern_utimesat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
2809            uap->tptr, UIO_USERSPACE));
2810 }
2811
2812 #ifndef _SYS_SYSPROTO_H_
2813 struct futimesat_args {
2814        int fd;
2815        const char * path;
2816        const struct timeval * times;
2817 };
2818 #endif
2819 int
2820 sys_futimesat(struct thread *td, struct futimesat_args *uap)
2821 {
2822
```

```
2823        return (kern_utimesat(td, uap->fd, uap->path, UIO_USERSPACE,
2824            uap->times, UIO_USERSPACE));
2825 }
2826
2827 int
2828 kern_utimesat(struct thread *td, int fd, char *path, enum uio_seg pathseg,
2829     struct timeval *tptr, enum uio_seg tptrseg)
2830 {
2831        struct nameidata nd;
2832        struct timespec ts[2];
2833        cap_rights_t rights;
2834        int error;
2835
2836        if ((error = getutimes(tptr, tptrseg, ts)) != 0)
2837                return (error);
2838        NDINIT_ATRIGHTS(&nd, LOOKUP, FOLLOW | AUDITVNODE1, pathseg, path, fd,
2839            cap_rights_init(&rights, CAP_FUTIMES), td);
2840
2841        if ((error = namei(&nd)) != 0)
2842                return (error);
2843        NDFREE(&nd, NDF_ONLY_PNBUF);
2844        error = setutimes(td, nd.ni_vp, ts, 2, tptr == NULL);
2845        vrele(nd.ni_vp);
2846        return (error);
2847 }
2848
2849 /*
2850  * Set the access and modification times of a file.
2851  */
2852 #ifndef _SYS_SYSPROTO_H_
2853 struct lutimes_args {
2854        char    *path;
2855        struct  timeval *tptr;
2856 };
2857 #endif
2858 int
2859 sys_lutimes(struct thread *td, struct lutimes_args *uap)
2860 {
2861
2862        return (kern_lutimes(td, uap->path, UIO_USERSPACE, uap->tptr,
2863            UIO_USERSPACE));
2864 }
2865
2866 int
2867 kern_lutimes(struct thread *td, char *path, enum uio_seg pathseg,
2868     struct timeval *tptr, enum uio_seg tptrseg)
2869 {
2870        struct timespec ts[2];
2871        struct nameidata nd;
2872        int error;
2873
2874        if ((error = getutimes(tptr, tptrseg, ts)) != 0)
2875                return (error);
2876        NDINIT(&nd, LOOKUP, NOFOLLOW | AUDITVNODE1, pathseg, path, td);
2877        if ((error = namei(&nd)) != 0)
2878                return (error);
2879        NDFREE(&nd, NDF_ONLY_PNBUF);
2880        error = setutimes(td, nd.ni_vp, ts, 2, tptr == NULL);
2881        vrele(nd.ni_vp);
2882        return (error);
2883 }
2884
2885 /*
2886  * Set the access and modification times of a file.
2887  */
2888 #ifndef _SYS_SYSPROTO_H_
2889 struct futimes_args {
2890        int     fd;
2891        struct  timeval *tptr;
2892 };
2893 #endif
2894 int
2895 sys_futimes(struct thread *td, struct futimes_args *uap)
2896 {
2897
2898        return (kern_futimes(td, uap->fd, uap->tptr, UIO_USERSPACE));
2899 }
2900
2901 int
2902 kern_futimes(struct thread *td, int fd, struct timeval *tptr,
2903     enum uio_seg tptrseg)
2904 {
2905        struct timespec ts[2];
2906        struct file *fp;
2907        cap_rights_t rights;
2908        int error;
2909
2910        AUDIT_ARG_FD(fd);
2911        error = getutimes(tptr, tptrseg, ts);
2912        if (error != 0)
2913                return (error);
2914        error = getvnode(td, fd, cap_rights_init(&rights, CAP_FUTIMES), &fp);
2915        if (error != 0)
2916                return (error);
2917 #ifdef AUDIT
2918        vn_lock(fp->f_vnode, LK_SHARED | LK_RETRY);
2919        AUDIT_ARG_VNODE1(fp->f_vnode);
2920        VOP_UNLOCK(fp->f_vnode, 0);
```

```
2921    #endif
2922            error = setutimes(td, fp->f_vnode, ts, 2, tptr == NULL);
2923            fdrop(fp, td);
2924            return (error);
2925    }
2926
2927    int
2928    sys_futimens(struct thread *td, struct futimens_args *uap)
2929    {
2930
2931            return (kern_futimens(td, uap->fd, uap->times, UIO_USERSPACE));
2932    }
2933
2934    int
2935    kern_futimens(struct thread *td, int fd, struct timespec *tptr,
2936        enum uio_seg tptrseg)
2937    {
2938            struct timespec ts[2];
2939            struct file *fp;
2940            cap_rights_t rights;
2941            int error, flags;
2942
2943            AUDIT_ARG_FD(fd);
2944            error = getutimens(tptr, tptrseg, ts, &flags);
2945            if (error != 0)
2946                    return (error);
2947            if (flags & UTIMENS_EXIT)
2948                    return (0);
2949            error = getvnode(td, fd, cap_rights_init(&rights, CAP_FUTIMES), &fp);
2950            if (error != 0)
2951                    return (error);
2952    #ifdef AUDIT
2953            vn_lock(fp->f_vnode, LK_SHARED | LK_RETRY);
2954            AUDIT_ARG_VNODE1(fp->f_vnode);
2955            VOP_UNLOCK(fp->f_vnode, 0);
2956    #endif
2957            error = setutimes(td, fp->f_vnode, ts, 2, flags & UTIMENS_NULL);
2958            fdrop(fp, td);
2959            return (error);
2960    }
2961
2962    int
2963    sys_utimensat(struct thread *td, struct utimensat_args *uap)
2964    {
2965
2966            return (kern_utimensat(td, uap->fd, uap->path, UIO_USERSPACE,
2967                uap->times, UIO_USERSPACE, uap->flag));
2968    }
2969
2970    int
2971    kern_utimensat(struct thread *td, int fd, char *path, enum uio_seg pathseg,
2972        struct timespec *tptr, enum uio_seg tptrseg, int flag)
2973    {
2974            struct nameidata nd;
2975            struct timespec ts[2];
2976            cap_rights_t rights;
2977            int error, flags;
2978
2979            if (flag & ~AT_SYMLINK_NOFOLLOW)
2980                    return (EINVAL);
2981
2982            if ((error = getutimens(tptr, tptrseg, ts, &flags)) != 0)
2983                    return (error);
2984            NDINIT_ATRIGHTS(&nd, LOOKUP, ((flag & AT_SYMLINK_NOFOLLOW) ? NOFOLLOW :
2985                FOLLOW) | AUDITVNODE1, pathseg, path, fd,
2986                cap_rights_init(&rights, CAP_FUTIMES), td);
2987            if ((error = namei(&nd)) != 0)
2988                    return (error);
2989            /*
2990             * We are allowed to call namei() regardless of 2xUTIME_OMIT.
2991             * POSIX states:
2992             * "If both tv_nsec fields are UTIME_OMIT... EACCESS may be detected."
2993             * "Search permission is denied by a component of the path prefix."
2994             */
2995            NDFREE(&nd, NDF_ONLY_PNBUF);
2996            if ((flags & UTIMENS_EXIT) == 0)
2997                    error = setutimes(td, nd.ni_vp, ts, 2, flags & UTIMENS_NULL);
2998            vrele(nd.ni_vp);
2999            return (error);
3000    }
3001
3002    /*
3003     * Truncate a file given its path name.
3004     */
3005    #ifndef _SYS_SYSPROTO_H_
3006    struct truncate_args {
3007            char    *path;
3008            int     pad;
3009            off_t   length;
3010    };
3011    #endif
3012    int
3013    sys_truncate(struct thread *td, struct truncate_args *uap)
3014    {
3015
3016            return (kern_truncate(td, uap->path, UIO_USERSPACE, uap->length));
3017    }
3018
```

```c
3019  int
3020  kern_truncate(struct thread *td, char *path, enum uio_seg pathseg, off_t length)
3021  {
3022          struct mount *mp;
3023          struct vnode *vp;
3024          void *rl_cookie;
3025          struct vattr vattr;
3026          struct nameidata nd;
3027          int error;
3028
3029          if (length < 0)
3030                  return(EINVAL);
3031          NDINIT(&nd, LOOKUP, FOLLOW | AUDITVNODE1, pathseg, path, td);
3032          if ((error = namei(&nd)) != 0)
3033                  return (error);
3034          vp = nd.ni_vp;
3035          rl_cookie = vn_rangelock_wlock(vp, 0, OFF_MAX);
3036          if ((error = vn_start_write(vp, &mp, V_WAIT | PCATCH)) != 0) {
3037                  vn_rangelock_unlock(vp, rl_cookie);
3038                  vrele(vp);
3039                  return (error);
3040          }
3041          NDFREE(&nd, NDF_ONLY_PNBUF);
3042          vn_lock(vp, LK_EXCLUSIVE | LK_RETRY);
3043          if (vp->v_type == VDIR)
3044                  error = EISDIR;
3045  #ifdef MAC
3046          else if ((error = mac_vnode_check_write(td->td_ucred, NOCRED, vp))) {
3047          }
3048  #endif
3049          else if ((error = vn_writechk(vp)) == 0 &&
3050              (error = VOP_ACCESS(vp, VWRITE, td->td_ucred, td)) == 0) {
3051                  VATTR_NULL(&vattr);
3052                  vattr.va_size = length;
3053                  error = VOP_SETATTR(vp, &vattr, td->td_ucred);
3054          }
3055          VOP_UNLOCK(vp, 0);
3056          vn_finished_write(mp);
3057          vn_rangelock_unlock(vp, rl_cookie);
3058          vrele(vp);
3059          return (error);
3060  }
3061
3062  #if defined(COMPAT_43)
3063  /*
3064   * Truncate a file given its path name.
3065   */
3066  #ifndef _SYS_SYSPROTO_H_
3067  struct otruncate_args {
3068          char    *path;
3069          long    length;
3070  };
3071  #endif
3072  int
3073  otruncate(struct thread *td, struct otruncate_args *uap)
3074  {
3075
3076          return (kern_truncate(td, uap->path, UIO_USERSPACE, uap->length));
3077  }
3078  #endif /* COMPAT_43 */
3079
3080  #if defined(COMPAT_FREEBSD6)
3081  /* Versions with the pad argument */
3082  int
3083  freebsd6_truncate(struct thread *td, struct freebsd6_truncate_args *uap)
3084  {
3085
3086          return (kern_truncate(td, uap->path, UIO_USERSPACE, uap->length));
3087  }
3088
3089  int
3090  freebsd6_ftruncate(struct thread *td, struct freebsd6_ftruncate_args *uap)
3091  {
3092
3093          return (kern_ftruncate(td, uap->fd, uap->length));
3094  }
3095  #endif
3096
3097  int
3098  kern_fsync(struct thread *td, int fd, bool fullsync)
3099  {
3100          struct vnode *vp;
3101          struct mount *mp;
3102          struct file *fp;
3103          cap_rights_t rights;
3104          int error, lock_flags;
3105
3106          AUDIT_ARG_FD(fd);
3107          error = getvnode(td, fd, cap_rights_init(&rights, CAP_FSYNC), &fp);
3108          if (error != 0)
3109                  return (error);
3110          vp = fp->f_vnode;
3111  #if 0
3112          if (!fullsync)
3113                  /* XXXKIB: compete outstanding aio writes */;
3114  #endif
3115          error = vn_start_write(vp, &mp, V_WAIT | PCATCH);
3116          if (error != 0)
```

```c
                goto drop;
        if (MNT_SHARED_WRITES(mp) ||
            ((mp == NULL) && MNT_SHARED_WRITES(vp->v_mount))) {
                lock_flags = LK_SHARED;
        } else {
                lock_flags = LK_EXCLUSIVE;
        }
        vn_lock(vp, lock_flags | LK_RETRY);
        AUDIT_ARG_VNODE1(vp);
        if (vp->v_object != NULL) {
                VM_OBJECT_WLOCK(vp->v_object);
                vm_object_page_clean(vp->v_object, 0, 0, 0);
                VM_OBJECT_WUNLOCK(vp->v_object);
        }
        error = fullsync ? VOP_FSYNC(vp, MNT_WAIT, td) : VOP_FDATASYNC(vp, td);
        VOP_UNLOCK(vp, 0);
        vn_finished_write(mp);
drop:
        fdrop(fp, td);
        return (error);
}

/*
 * Sync an open file.
 */
#ifndef _SYS_SYSPROTO_H_
struct fsync_args {
        int     fd;
};
#endif
int
sys_fsync(struct thread *td, struct fsync_args *uap)
{

        return (kern_fsync(td, uap->fd, true));
}

int
sys_fdatasync(struct thread *td, struct fdatasync_args *uap)
{

        return (kern_fsync(td, uap->fd, false));
}

/*
 * Rename files.  Source and destination must either both be directories, or
 * both not be directories.  If target is a directory, it must be empty.
 */
#ifndef _SYS_SYSPROTO_H_
struct rename_args {
        char    *from;
        char    *to;
};
#endif
int
sys_rename(struct thread *td, struct rename_args *uap)
{

        return (kern_renameat(td, AT_FDCWD, uap->from, AT_FDCWD,
            uap->to, UIO_USERSPACE));
}

#ifndef _SYS_SYSPROTO_H_
struct renameat_args {
        int     oldfd;
        char    *old;
        int     newfd;
        char    *new;
};
#endif
int
sys_renameat(struct thread *td, struct renameat_args *uap)
{

        return (kern_renameat(td, uap->oldfd, uap->old, uap->newfd, uap->new,
            UIO_USERSPACE));
}

int
kern_renameat(struct thread *td, int oldfd, char *old, int newfd, char *new,
    enum uio_seg pathseg)
{
        struct mount *mp = NULL;
        struct vnode *tvp, *fvp, *tdvp;
        struct nameidata fromnd, tond;
        cap_rights_t rights;
        int error;

again:
        bwillwrite();
#ifdef MAC
        NDINIT_ATRIGHTS(&fromnd, DELETE, LOCKPARENT | LOCKLEAF | SAVESTART |
            AUDITVNODE1, pathseg, old, oldfd,
            cap_rights_init(&rights, CAP_RENAMEAT_SOURCE), td);
#else
        NDINIT_ATRIGHTS(&fromnd, DELETE, WANTPARENT | SAVESTART | AUDITVNODE1,
            pathseg, old, oldfd,
            cap_rights_init(&rights, CAP_RENAMEAT_SOURCE), td);
```

```c
3215  #endif
3216
3217          if ((error = namei(&fromnd)) != 0)
3218                  return (error);
3219  #ifdef MAC
3220          error = mac_vnode_check_rename_from(td->td_ucred, fromnd.ni_dvp,
3221              fromnd.ni_vp, &fromnd.ni_cnd);
3222          VOP_UNLOCK(fromnd.ni_dvp, 0);
3223          if (fromnd.ni_dvp != fromnd.ni_vp)
3224                  VOP_UNLOCK(fromnd.ni_vp, 0);
3225  #endif
3226          fvp = fromnd.ni_vp;
3227          NDINIT_ATRIGHTS(&tond, RENAME, LOCKPARENT | LOCKLEAF | NOCACHE |
3228              SAVESTART | AUDITVNODE2, pathseg, new, newfd,
3229              cap_rights_init(&rights, CAP_RENAMEAT_TARGET), td);
3230          if (fromnd.ni_vp->v_type == VDIR)
3231                  tond.ni_cnd.cn_flags |= WILLBEDIR;
3232          if ((error = namei(&tond)) != 0) {
3233                  /* Translate error code for rename("dir1", "dir2/."). */
3234                  if (error == EISDIR && fvp->v_type == VDIR)
3235                          error = EINVAL;
3236                  NDFREE(&fromnd, NDF_ONLY_PNBUF);
3237                  vrele(fromnd.ni_dvp);
3238                  vrele(fvp);
3239                  goto out1;
3240          }
3241          tdvp = tond.ni_dvp;
3242          tvp = tond.ni_vp;
3243          error = vn_start_write(fvp, &mp, V_NOWAIT);
3244          if (error != 0) {
3245                  NDFREE(&fromnd, NDF_ONLY_PNBUF);
3246                  NDFREE(&tond, NDF_ONLY_PNBUF);
3247                  if (tvp != NULL)
3248                          vput(tvp);
3249                  if (tdvp == tvp)
3250                          vrele(tdvp);
3251                  else
3252                          vput(tdvp);
3253                  vrele(fromnd.ni_dvp);
3254                  vrele(fvp);
3255                  vrele(tond.ni_startdir);
3256                  if (fromnd.ni_startdir != NULL)
3257                          vrele(fromnd.ni_startdir);
3258                  error = vn_start_write(NULL, &mp, V_XSLEEP | PCATCH);
3259                  if (error != 0)
3260                          return (error);
3261                  goto again;
3262          }
3263          if (tvp != NULL) {
3264                  if (fvp->v_type == VDIR && tvp->v_type != VDIR) {
3265                          error = ENOTDIR;
3266                          goto out;
3267                  } else if (fvp->v_type != VDIR && tvp->v_type == VDIR) {
3268                          error = EISDIR;
3269                          goto out;
3270                  }
3271  #ifdef CAPABILITIES
3272                  if (newfd != AT_FDCWD) {
3273                          /*
3274                           * If the target already exists we require CAP_UNLINKAT
3275                           * from 'newfd'.
3276                           */
3277                          error = cap_check(&tond.ni_filecaps.fc_rights,
3278                              cap_rights_init(&rights, CAP_UNLINKAT));
3279                          if (error != 0)
3280                                  goto out;
3281                  }
3282  #endif
3283          }
3284          if (fvp == tdvp) {
3285                  error = EINVAL;
3286                  goto out;
3287          }
3288          /*
3289           * If the source is the same as the destination (that is, if they
3290           * are links to the same vnode), then there is nothing to do.
3291           */
3292          if (fvp == tvp)
3293                  error = -1;
3294  #ifdef MAC
3295          else
3296                  error = mac_vnode_check_rename_to(td->td_ucred, tdvp,
3297                      tond.ni_vp, fromnd.ni_dvp == tdvp, &tond.ni_cnd);
3298  #endif
3299  out:
3300          if (error == 0) {
3301                  error = VOP_RENAME(fromnd.ni_dvp, fromnd.ni_vp, &fromnd.ni_cnd,
3302                      tond.ni_dvp, tond.ni_vp, &tond.ni_cnd);
3303                  NDFREE(&fromnd, NDF_ONLY_PNBUF);
3304                  NDFREE(&tond, NDF_ONLY_PNBUF);
3305          } else {
3306                  NDFREE(&fromnd, NDF_ONLY_PNBUF);
3307                  NDFREE(&tond, NDF_ONLY_PNBUF);
3308                  if (tvp != NULL)
3309                          vput(tvp);
3310                  if (tdvp == tvp)
3311                          vrele(tdvp);
3312                  else
```

```c
                                vput(tdvp);
                        vrele(fromnd.ni_dvp);
                        vrele(fvp);
                }
                vrele(tond.ni_startdir);
                vn_finished_write(mp);
out1:
                if (fromnd.ni_startdir)
                        vrele(fromnd.ni_startdir);
                if (error == -1)
                        return (0);
                return (error);
}

/*
 * Make a directory file.
 */
#ifndef _SYS_SYSPROTO_H_
struct mkdir_args {
        char    *path;
        int     mode;
};
#endif
int
sys_mkdir(struct thread *td, struct mkdir_args *uap)
{

                return (kern_mkdirat(td, AT_FDCWD, uap->path, UIO_USERSPACE,
                    uap->mode));
}

#ifndef _SYS_SYSPROTO_H_
struct mkdirat_args {
        int     fd;
        char    *path;
        mode_t  mode;
};
#endif
int
sys_mkdirat(struct thread *td, struct mkdirat_args *uap)
{

                return (kern_mkdirat(td, uap->fd, uap->path, UIO_USERSPACE, uap->mode));
}

int
kern_mkdirat(struct thread *td, int fd, char *path, enum uio_seg segflg,
    int mode)
{
                struct mount *mp;
                struct vnode *vp;
                struct vattr vattr;
                struct nameidata nd;
                cap_rights_t rights;
                int error;

                AUDIT_ARG_MODE(mode);
restart:
                bwillwrite();
                NDINIT_ATRIGHTS(&nd, CREATE, LOCKPARENT | SAVENAME | AUDITVNODE1 |
                    NOCACHE, segflg, path, fd, cap_rights_init(&rights, CAP_MKDIRAT),
                    td);
                nd.ni_cnd.cn_flags |= WILLBEDIR;
                if ((error = namei(&nd)) != 0)
                        return (error);
                vp = nd.ni_vp;
                if (vp != NULL) {
                        NDFREE(&nd, NDF_ONLY_PNBUF);
                        /*
                         * XXX namei called with LOCKPARENT but not LOCKLEAF has
                         * the strange behaviour of leaving the vnode unlocked
                         * if the target is the same vnode as the parent.
                         */
                        if (vp == nd.ni_dvp)
                                vrele(nd.ni_dvp);
                        else
                                vput(nd.ni_dvp);
                        vrele(vp);
                        return (EEXIST);
                }
                if (vn_start_write(nd.ni_dvp, &mp, V_NOWAIT) != 0) {
                        NDFREE(&nd, NDF_ONLY_PNBUF);
                        vput(nd.ni_dvp);
                        if ((error = vn_start_write(NULL, &mp, V_XSLEEP | PCATCH)) != 0)
                                return (error);
                        goto restart;
                }
                VATTR_NULL(&vattr);
                vattr.va_type = VDIR;
                vattr.va_mode = (mode & ACCESSPERMS) &~ td->td_proc->p_fd->fd_cmask;
#ifdef MAC
                error = mac_vnode_check_create(td->td_ucred, nd.ni_dvp, &nd.ni_cnd,
                    &vattr);
                if (error != 0)
                        goto out;
#endif
                error = VOP_MKDIR(nd.ni_dvp, &nd.ni_vp, &nd.ni_cnd, &vattr);
#ifdef MAC
```

```
3411    out:
3412    #endif
3413            NDFREE(&nd, NDF_ONLY_PNBUF);
3414            vput(nd.ni_dvp);
3415            if (error == 0)
3416                    vput(nd.ni_vp);
3417            vn_finished_write(mp);
3418            return (error);
3419    }
3420
3421    /*
3422     * Remove a directory file.
3423     */
3424    #ifndef _SYS_SYSPROTO_H_
3425    struct rmdir_args {
3426            char    *path;
3427    };
3428    #endif
3429    int
3430    sys_rmdir(struct thread *td, struct rmdir_args *uap)
3431    {
3432
3433            return (kern_rmdirat(td, AT_FDCWD, uap->path, UIO_USERSPACE));
3434    }
3435
3436    int
3437    kern_rmdirat(struct thread *td, int fd, char *path, enum uio_seg pathseg)
3438    {
3439            struct mount *mp;
3440            struct vnode *vp;
3441            struct nameidata nd;
3442            cap_rights_t rights;
3443            int error;
3444
3445    restart:
3446            bwillwrite();
3447            NDINIT_ATRIGHTS(&nd, DELETE, LOCKPARENT | LOCKLEAF | AUDITVNODE1,
3448                pathseg, path, fd, cap_rights_init(&rights, CAP_UNLINKAT), td);
3449            if ((error = namei(&nd)) != 0)
3450                    return (error);
3451            vp = nd.ni_vp;
3452            if (vp->v_type != VDIR) {
3453                    error = ENOTDIR;
3454                    goto out;
3455            }
3456            /*
3457             * No rmdir "." please.
3458             */
3459            if (nd.ni_dvp == vp) {
3460                    error = EINVAL;
3461                    goto out;
3462            }
3463            /*
3464             * The root of a mounted filesystem cannot be deleted.
3465             */
3466            if (vp->v_vflag & VV_ROOT) {
3467                    error = EBUSY;
3468                    goto out;
3469            }
3470    #ifdef MAC
3471            error = mac_vnode_check_unlink(td->td_ucred, nd.ni_dvp, vp,
3472                &nd.ni_cnd);
3473            if (error != 0)
3474                    goto out;
3475    #endif
3476            if (vn_start_write(nd.ni_dvp, &mp, V_NOWAIT) != 0) {
3477                    NDFREE(&nd, NDF_ONLY_PNBUF);
3478                    vput(vp);
3479                    if (nd.ni_dvp == vp)
3480                            vrele(nd.ni_dvp);
3481                    else
3482                            vput(nd.ni_dvp);
3483                    if ((error = vn_start_write(NULL, &mp, V_XSLEEP | PCATCH)) != 0)
3484                            return (error);
3485                    goto restart;
3486            }
3487            vfs_notify_upper(vp, VFS_NOTIFY_UPPER_UNLINK);
3488            error = VOP_RMDIR(nd.ni_dvp, nd.ni_vp, &nd.ni_cnd);
3489            vn_finished_write(mp);
3490    out:
3491            NDFREE(&nd, NDF_ONLY_PNBUF);
3492            vput(vp);
3493            if (nd.ni_dvp == vp)
3494                    vrele(nd.ni_dvp);
3495            else
3496                    vput(nd.ni_dvp);
3497            return (error);
3498    }
3499
3500    #ifdef COMPAT_43
3501    /*
3502     * Read a block of directory entries in a filesystem independent format.
3503     */
3504    #ifndef _SYS_SYSPROTO_H_
3505    struct ogetdirentries_args {
3506            int     fd;
3507            char    *buf;
3508            u_int   count;
```

```
3509          long     *basep;
3510   };
3511   #endif
3512   int
3513   ogetdirentries(struct thread *td, struct ogetdirentries_args *uap)
3514   {
3515          long loff;
3516          int error;
3517
3518          error = kern_ogetdirentries(td, uap, &loff);
3519          if (error == 0)
3520                  error = copyout(&loff, uap->basep, sizeof(long));
3521          return (error);
3522   }
3523
3524   int
3525   kern_ogetdirentries(struct thread *td, struct ogetdirentries_args *uap,
3526       long *ploff)
3527   {
3528          struct vnode *vp;
3529          struct file *fp;
3530          struct uio auio, kuio;
3531          struct iovec aiov, kiov;
3532          struct dirent *dp, *edp;
3533          cap_rights_t rights;
3534          caddr_t dirbuf;
3535          int error, eofflag, readcnt;
3536          long loff;
3537          off_t foffset;
3538
3539          /* XXX arbitrary sanity limit on `count'. */
3540          if (uap->count > 64 * 1024)
3541                  return (EINVAL);
3542          error = getvnode(td, uap->fd, cap_rights_init(&rights, CAP_READ), &fp);
3543          if (error != 0)
3544                  return (error);
3545          if ((fp->f_flag & FREAD) == 0) {
3546                  fdrop(fp, td);
3547                  return (EBADF);
3548          }
3549          vp = fp->f_vnode;
3550          foffset = foffset_lock(fp, 0);
3551   unionread:
3552          if (vp->v_type != VDIR) {
3553                  foffset_unlock(fp, foffset, 0);
3554                  fdrop(fp, td);
3555                  return (EINVAL);
3556          }
3557          aiov.iov_base = uap->buf;
3558          aiov.iov_len = uap->count;
3559          auio.uio_iov = &aiov;
3560          auio.uio_iovcnt = 1;
3561          auio.uio_rw = UIO_READ;
3562          auio.uio_segflg = UIO_USERSPACE;
3563          auio.uio_td = td;
3564          auio.uio_resid = uap->count;
3565          vn_lock(vp, LK_SHARED | LK_RETRY);
3566          loff = auio.uio_offset = foffset;
3567   #ifdef MAC
3568          error = mac_vnode_check_readdir(td->td_ucred, vp);
3569          if (error != 0) {
3570                  VOP_UNLOCK(vp, 0);
3571                  foffset_unlock(fp, foffset, FOF_NOUPDATE);
3572                  fdrop(fp, td);
3573                  return (error);
3574          }
3575   #endif
3576   #      if (BYTE_ORDER != LITTLE_ENDIAN)
3577                  if (vp->v_mount->mnt_maxsymlinklen <= 0) {
3578                          error = VOP_READDIR(vp, &auio, fp->f_cred, &eofflag,
3579                              NULL, NULL);
3580                          foffset = auio.uio_offset;
3581                  } else
3582   #      endif
3583          {
3584                  kuio = auio;
3585                  kuio.uio_iov = &kiov;
3586                  kuio.uio_segflg = UIO_SYSSPACE;
3587                  kiov.iov_len = uap->count;
3588                  dirbuf = malloc(uap->count, M_TEMP, M_WAITOK);
3589                  kiov.iov_base = dirbuf;
3590                  error = VOP_READDIR(vp, &kuio, fp->f_cred, &eofflag,
3591                              NULL, NULL);
3592                  foffset = kuio.uio_offset;
3593                  if (error == 0) {
3594                          readcnt = uap->count - kuio.uio_resid;
3595                          edp = (struct dirent *)&dirbuf[readcnt];
3596                          for (dp = (struct dirent *)dirbuf; dp < edp; ) {
3597   #                              if (BYTE_ORDER == LITTLE_ENDIAN)
3598                                          /*
3599                                           * The expected low byte of
3600                                           * dp->d_namlen is our dp->d_type.
3601                                           * The high MBZ byte of dp->d_namlen
3602                                           * is our dp->d_namlen.
3603                                           */
3604                                          dp->d_type = dp->d_namlen;
3605                                          dp->d_namlen = 0;
3606   #                              else
```

```
3607                                         /*
3608                                          * The dp->d_type is the high byte
3609                                          * of the expected dp->d_namlen,
3610                                          * so must be zero'ed.
3611                                          */
3612                                         dp->d_type = 0;
3613 #                                endif
3614                                 if (dp->d_reclen > 0) {
3615                                         dp = (struct dirent *)
3616                                             ((char *)dp + dp->d_reclen);
3617                                 } else {
3618                                         error = EIO;
3619                                         break;
3620                                 }
3621                         }
3622                         if (dp >= edp)
3623                                 error = uiomove(dirbuf, readcnt, &auio);
3624                 }
3625                 free(dirbuf, M_TEMP);
3626         }
3627         if (error != 0) {
3628                 VOP_UNLOCK(vp, 0);
3629                 foffset_unlock(fp, foffset, 0);
3630                 fdrop(fp, td);
3631                 return (error);
3632         }
3633         if (uap->count == auio.uio_resid &&
3634             (vp->v_vflag & VV_ROOT) &&
3635             (vp->v_mount->mnt_flag & MNT_UNION)) {
3636                 struct vnode *tvp = vp;
3637                 vp = vp->v_mount->mnt_vnodecovered;
3638                 VREF(vp);
3639                 fp->f_vnode = vp;
3640                 fp->f_data = vp;
3641                 foffset = 0;
3642                 vput(tvp);
3643                 goto unionread;
3644         }
3645         VOP_UNLOCK(vp, 0);
3646         foffset_unlock(fp, foffset, 0);
3647         fdrop(fp, td);
3648         td->td_retval[0] = uap->count - auio.uio_resid;
3649         if (error == 0)
3650                 *ploff = loff;
3651         return (error);
3652 }
3653 #endif /* COMPAT_43 */
3654
3655 /*
3656  * Read a block of directory entries in a filesystem independent format.
3657  */
3658 #ifndef _SYS_SYSPROTO_H_
3659 struct getdirentries_args {
3660         int     fd;
3661         char    *buf;
3662         u_int   count;
3663         long    *basep;
3664 };
3665 #endif
3666 int
3667 sys_getdirentries(struct thread *td, struct getdirentries_args *uap)
3668 {
3669         long base;
3670         int error;
3671
3672         error = kern_getdirentries(td, uap->fd, uap->buf, uap->count, &base,
3673             NULL, UIO_USERSPACE);
3674         if (error != 0)
3675                 return (error);
3676         if (uap->basep != NULL)
3677                 error = copyout(&base, uap->basep, sizeof(long));
3678         return (error);
3679 }
3680
3681 int
3682 kern_getdirentries(struct thread *td, int fd, char *buf, u_int count,
3683     long *basep, ssize_t *residp, enum uio_seg bufseg)
3684 {
3685         struct vnode *vp;
3686         struct file *fp;
3687         struct uio auio;
3688         struct iovec aiov;
3689         cap_rights_t rights;
3690         long loff;
3691         int error, eofflag;
3692         off_t foffset;
3693
3694         AUDIT_ARG_FD(fd);
3695         if (count > IOSIZE_MAX)
3696                 return (EINVAL);
3697         auio.uio_resid = count;
3698         error = getvnode(td, fd, cap_rights_init(&rights, CAP_READ), &fp);
3699         if (error != 0)
3700                 return (error);
3701         if ((fp->f_flag & FREAD) == 0) {
3702                 fdrop(fp, td);
3703                 return (EBADF);
3704         }
```

```
3705            vp = fp->f_vnode;
3706            foffset = foffset_lock(fp, 0);
3707    unionread:
3708            if (vp->v_type != VDIR) {
3709                    error = EINVAL;
3710                    goto fail;
3711            }
3712            aiov.iov_base = buf;
3713            aiov.iov_len = count;
3714            auio.uio_iov = &aiov;
3715            auio.uio_iovcnt = 1;
3716            auio.uio_rw = UIO_READ;
3717            auio.uio_segflg = bufseg;
3718            auio.uio_td = td;
3719            vn_lock(vp, LK_SHARED | LK_RETRY);
3720            AUDIT_ARG_VNODE1(vp);
3721            loff = auio.uio_offset = foffset;
3722    #ifdef MAC
3723            error = mac_vnode_check_readdir(td->td_ucred, vp);
3724            if (error == 0)
3725    #endif
3726                    error = VOP_READDIR(vp, &auio, fp->f_cred, &eofflag, NULL,
3727                        NULL);
3728            foffset = auio.uio_offset;
3729            if (error != 0) {
3730                    VOP_UNLOCK(vp, 0);
3731                    goto fail;
3732            }
3733            if (count == auio.uio_resid &&
3734                (vp->v_vflag & VV_ROOT) &&
3735                (vp->v_mount->mnt_flag & MNT_UNION)) {
3736                    struct vnode *tvp = vp;
3737
3738                    vp = vp->v_mount->mnt_vnodecovered;
3739                    VREF(vp);
3740                    fp->f_vnode = vp;
3741                    fp->f_data = vp;
3742                    foffset = 0;
3743                    vput(tvp);
3744                    goto unionread;
3745            }
3746            VOP_UNLOCK(vp, 0);
3747            *basep = loff;
3748            if (residp != NULL)
3749                    *residp = auio.uio_resid;
3750            td->td_retval[0] = count - auio.uio_resid;
3751    fail:
3752            foffset_unlock(fp, foffset, 0);
3753            fdrop(fp, td);
3754            return (error);
3755    }
3756
3757    #ifndef _SYS_SYSPROTO_H_
3758    struct getdents_args {
3759            int fd;
3760            char *buf;
3761            size_t count;
3762    };
3763    #endif
3764    int
3765    sys_getdents(struct thread *td, struct getdents_args *uap)
3766    {
3767            struct getdirentries_args ap;
3768
3769            ap.fd = uap->fd;
3770            ap.buf = uap->buf;
3771            ap.count = uap->count;
3772            ap.basep = NULL;
3773            return (sys_getdirentries(td, &ap));
3774    }
3775
3776    /*
3777     * Set the mode mask for creation of filesystem nodes.
3778     */
3779    #ifndef _SYS_SYSPROTO_H_
3780    struct umask_args {
3781            int     newmask;
3782    };
3783    #endif
3784    int
3785    sys_umask(struct thread *td, struct umask_args *uap)
3786    {
3787            struct filedesc *fdp;
3788
3789            fdp = td->td_proc->p_fd;
3790            FILEDESC_XLOCK(fdp);
3791            td->td_retval[0] = fdp->fd_cmask;
3792            fdp->fd_cmask = uap->newmask & ALLPERMS;
3793            FILEDESC_XUNLOCK(fdp);
3794            return (0);
3795    }
3796
3797    /*
3798     * Void all references to file by ripping underlying filesystem away from
3799     * vnode.
3800     */
3801    #ifndef _SYS_SYSPROTO_H_
3802    struct revoke_args {
```

```
3803            char    *path;
3804    };
3805    #endif
3806    int
3807    sys_revoke(struct thread *td, struct revoke_args *uap)
3808    {
3809            struct vnode *vp;
3810            struct vattr vattr;
3811            struct nameidata nd;
3812            int error;
3813
3814            NDINIT(&nd, LOOKUP, FOLLOW | LOCKLEAF | AUDITVNODE1, UIO_USERSPACE,
3815                uap->path, td);
3816            if ((error = namei(&nd)) != 0)
3817                    return (error);
3818            vp = nd.ni_vp;
3819            NDFREE(&nd, NDF_ONLY_PNBUF);
3820            if (vp->v_type != VCHR || vp->v_rdev == NULL) {
3821                    error = EINVAL;
3822                    goto out;
3823            }
3824    #ifdef MAC
3825            error = mac_vnode_check_revoke(td->td_ucred, vp);
3826            if (error != 0)
3827                    goto out;
3828    #endif
3829            error = VOP_GETATTR(vp, &vattr, td->td_ucred);
3830            if (error != 0)
3831                    goto out;
3832            if (td->td_ucred->cr_uid != vattr.va_uid) {
3833                    error = priv_check(td, PRIV_VFS_ADMIN);
3834                    if (error != 0)
3835                            goto out;
3836            }
3837            if (vcount(vp) > 1)
3838                    VOP_REVOKE(vp, REVOKEALL);
3839    out:
3840            vput(vp);
3841            return (error);
3842    }
3843
3844    /*
3845     * Convert a user file descriptor to a kernel file entry and check that, if it
3846     * is a capability, the correct rights are present. A reference on the file
3847     * entry is held upon returning.
3848     */
3849    int
3850    getvnode(struct thread *td, int fd, cap_rights_t *rightsp, struct file **fpp)
3851    {
3852            struct file *fp;
3853            int error;
3854
3855            error = fget_unlocked(td->td_proc->p_fd, fd, rightsp, &fp, NULL);
3856            if (error != 0)
3857                    return (error);
3858
3859            /*
3860             * The file could be not of the vnode type, or it may be not
3861             * yet fully initialized, in which case the f_vnode pointer
3862             * may be set, but f_ops is still badfileops.  E.g.,
3863             * devfs_open() transiently create such situation to
3864             * facilitate csw d_fdopen().
3865             *
3866             * Dupfdopen() handling in kern_openat() installs the
3867             * half-baked file into the process descriptor table, allowing
3868             * other thread to dereference it. Guard against the race by
3869             * checking f_ops.
3870             */
3871            if (fp->f_vnode == NULL || fp->f_ops == &badfileops) {
3872                    fdrop(fp, td);
3873                    return (EINVAL);
3874            }
3875            *fpp = fp;
3876            return (0);
3877    }
3878
3879
3880    /*
3881     * Get an (NFS) file handle.
3882     */
3883    #ifndef _SYS_SYSPROTO_H_
3884    struct lgetfh_args {
3885            char    *fname;
3886            fhandle_t *fhp;
3887    };
3888    #endif
3889    int
3890    sys_lgetfh(struct thread *td, struct lgetfh_args *uap)
3891    {
3892            struct nameidata nd;
3893            fhandle_t fh;
3894            struct vnode *vp;
3895            int error;
3896
3897            error = priv_check(td, PRIV_VFS_GETFH);
3898            if (error != 0)
3899                    return (error);
3900            NDINIT(&nd, LOOKUP, NOFOLLOW | LOCKLEAF | AUDITVNODE1, UIO_USERSPACE,
```

```
3901             uap->fname, td);
3902         error = namei(&nd);
3903         if (error != 0)
3904             return (error);
3905         NDFREE(&nd, NDF_ONLY_PNBUF);
3906         vp = nd.ni_vp;
3907         bzero(&fh, sizeof(fh));
3908         fh.fh_fsid = vp->v_mount->mnt_stat.f_fsid;
3909         error = VOP_VPTOFH(vp, &fh.fh_fid);
3910         vput(vp);
3911         if (error == 0)
3912             error = copyout(&fh, uap->fhp, sizeof (fh));
3913         return (error);
3914     }
3915
3916     #ifndef _SYS_SYSPROTO_H_
3917     struct getfh_args {
3918         char    *fname;
3919         fhandle_t *fhp;
3920     };
3921     #endif
3922     int
3923     sys_getfh(struct thread *td, struct getfh_args *uap)
3924     {
3925         struct nameidata nd;
3926         fhandle_t fh;
3927         struct vnode *vp;
3928         int error;
3929
3930         error = priv_check(td, PRIV_VFS_GETFH);
3931         if (error != 0)
3932             return (error);
3933         NDINIT(&nd, LOOKUP, FOLLOW | LOCKLEAF | AUDITVNODE1, UIO_USERSPACE,
3934             uap->fname, td);
3935         error = namei(&nd);
3936         if (error != 0)
3937             return (error);
3938         NDFREE(&nd, NDF_ONLY_PNBUF);
3939         vp = nd.ni_vp;
3940         bzero(&fh, sizeof(fh));
3941         fh.fh_fsid = vp->v_mount->mnt_stat.f_fsid;
3942         error = VOP_VPTOFH(vp, &fh.fh_fid);
3943         vput(vp);
3944         if (error == 0)
3945             error = copyout(&fh, uap->fhp, sizeof (fh));
3946         return (error);
3947     }
3948
3949     /*
3950      * syscall for the rpc.lockd to use to translate a NFS file handle into an
3951      * open descriptor.
3952      *
3953      * warning: do not remove the priv_check() call or this becomes one giant
3954      * security hole.
3955      */
3956     #ifndef _SYS_SYSPROTO_H_
3957     struct fhopen_args {
3958         const struct fhandle *u_fhp;
3959         int flags;
3960     };
3961     #endif
3962     int
3963     sys_fhopen(struct thread *td, struct fhopen_args *uap)
3964     {
3965         struct mount *mp;
3966         struct vnode *vp;
3967         struct fhandle fhp;
3968         struct file *fp;
3969         int fmode, error;
3970         int indx;
3971
3972         error = priv_check(td, PRIV_VFS_FHOPEN);
3973         if (error != 0)
3974             return (error);
3975         indx = -1;
3976         fmode = FFLAGS(uap->flags);
3977         /* why not allow a non-read/write open for our lockd? */
3978         if (((fmode & (FREAD | FWRITE)) == 0) || (fmode & O_CREAT))
3979             return (EINVAL);
3980         error = copyin(uap->u_fhp, &fhp, sizeof(fhp));
3981         if (error != 0)
3982             return(error);
3983         /* find the mount point */
3984         mp = vfs_busyfs(&fhp.fh_fsid);
3985         if (mp == NULL)
3986             return (ESTALE);
3987         /* now give me my vnode, it gets returned to me locked */
3988         error = VFS_FHTOVP(mp, &fhp.fh_fid, LK_EXCLUSIVE, &vp);
3989         vfs_unbusy(mp);
3990         if (error != 0)
3991             return (error);
3992
3993         error = falloc_noinstall(td, &fp);
3994         if (error != 0) {
3995             vput(vp);
3996             return (error);
3997         }
3998         /*
```

```
3999            * An extra reference on `fp' has been held for us by
4000            * falloc_noinstall().
4001            */

4003  #ifdef INVARIANTS
4004           td->td_dupfd = -1;
4005  #endif
4006           error = vn_open_vnode(vp, fmode, td->td_ucred, td, fp);
4007           if (error != 0) {
4008                   KASSERT(fp->f_ops == &badfileops,
4009                       ("VOP_OPEN in fhopen() set f_ops"));
4010                   KASSERT(td->td_dupfd < 0,
4011                       ("fhopen() encountered fdopen()"));

4013                   vput(vp);
4014                   goto bad;
4015           }
4016  #ifdef INVARIANTS
4017           td->td_dupfd = 0;
4018  #endif
4019           fp->f_vnode = vp;
4020           fp->f_seqcount = 1;
4021           finit(fp, (fmode & FMASK) | (fp->f_flag & FHASLOCK), DTYPE_VNODE, vp,
4022               &vnops);
4023           VOP_UNLOCK(vp, 0);
4024           if ((fmode & O_TRUNC) != 0) {
4025                   error = fo_truncate(fp, 0, td->td_ucred, td);
4026                   if (error != 0)
4027                           goto bad;
4028           }

4030           error = finstall(td, fp, &indx, fmode, NULL);
4031  bad:
4032           fdrop(fp, td);
4033           td->td_retval[0] = indx;
4034           return (error);
4035  }

4037  /*
4038   * Stat an (NFS) file handle.
4039   */
4040  #ifndef _SYS_SYSPROTO_H_
4041  struct fhstat_args {
4042           struct fhandle *u_fhp;
4043           struct stat *sb;
4044  };
4045  #endif
4046  int
4047  sys_fhstat(struct thread *td, struct fhstat_args *uap)
4048  {
4049           struct stat sb;
4050           struct fhandle fh;
4051           int error;

4053           error = copyin(uap->u_fhp, &fh, sizeof(fh));
4054           if (error != 0)
4055                   return (error);
4056           error = kern_fhstat(td, fh, &sb);
4057           if (error == 0)
4058                   error = copyout(&sb, uap->sb, sizeof(sb));
4059           return (error);
4060  }

4062  int
4063  kern_fhstat(struct thread *td, struct fhandle fh, struct stat *sb)
4064  {
4065           struct mount *mp;
4066           struct vnode *vp;
4067           int error;

4069           error = priv_check(td, PRIV_VFS_FHSTAT);
4070           if (error != 0)
4071                   return (error);
4072           if ((mp = vfs_busyfs(&fh.fh_fsid)) == NULL)
4073                   return (ESTALE);
4074           error = VFS_FHTOVP(mp, &fh.fh_fid, LK_EXCLUSIVE, &vp);
4075           vfs_unbusy(mp);
4076           if (error != 0)
4077                   return (error);
4078           error = vn_stat(vp, sb, td->td_ucred, NOCRED, td);
4079           vput(vp);
4080           return (error);
4081  }

4083  /*
4084   * Implement fstatfs() for (NFS) file handles.
4085   */
4086  #ifndef _SYS_SYSPROTO_H_
4087  struct fhstatfs_args {
4088           struct fhandle *u_fhp;
4089           struct statfs *buf;
4090  };
4091  #endif
4092  int
4093  sys_fhstatfs(struct thread *td, struct fhstatfs_args *uap)
4094  {
4095           struct statfs *sfp;
4096           fhandle_t fh;
```

```
4097            int error;
4098
4099            error = copyin(uap->u_fhp, &fh, sizeof(fhandle_t));
4100            if (error != 0)
4101                    return (error);
4102            sfp = malloc(sizeof(struct statfs), M_STATFS, M_WAITOK);
4103            error = kern_fhstatfs(td, fh, sfp);
4104            if (error == 0)
4105                    error = copyout(sfp, uap->buf, sizeof(*sfp));
4106            free(sfp, M_STATFS);
4107            return (error);
4108    }
4109
4110    int
4111    kern_fhstatfs(struct thread *td, fhandle_t fh, struct statfs *buf)
4112    {
4113            struct statfs *sp;
4114            struct mount *mp;
4115            struct vnode *vp;
4116            int error;
4117
4118            error = priv_check(td, PRIV_VFS_FHSTATFS);
4119            if (error != 0)
4120                    return (error);
4121            if ((mp = vfs_busyfs(&fh.fh_fsid)) == NULL)
4122                    return (ESTALE);
4123            error = VFS_FHTOVP(mp, &fh.fh_fid, LK_EXCLUSIVE, &vp);
4124            if (error != 0) {
4125                    vfs_unbusy(mp);
4126                    return (error);
4127            }
4128            vput(vp);
4129            error = prison_canseemount(td->td_ucred, mp);
4130            if (error != 0)
4131                    goto out;
4132    #ifdef MAC
4133            error = mac_mount_check_stat(td->td_ucred, mp);
4134            if (error != 0)
4135                    goto out;
4136    #endif
4137            /*
4138             * Set these in case the underlying filesystem fails to do so.
4139             */
4140            sp = &mp->mnt_stat;
4141            sp->f_version = STATFS_VERSION;
4142            sp->f_namemax = NAME_MAX;
4143            sp->f_flags = mp->mnt_flag & MNT_VISFLAGMASK;
4144            error = VFS_STATFS(mp, sp);
4145            if (error == 0)
4146                    *buf = *sp;
4147    out:
4148            vfs_unbusy(mp);
4149            return (error);
4150    }
4151
4152    int
4153    kern_posix_fallocate(struct thread *td, int fd, off_t offset, off_t len)
4154    {
4155            struct file *fp;
4156            struct mount *mp;
4157            struct vnode *vp;
4158            cap_rights_t rights;
4159            off_t olen, ooffset;
4160            int error;
4161
4162            if (offset < 0 || len <= 0)
4163                    return (EINVAL);
4164            /* Check for wrap. */
4165            if (offset > OFF_MAX - len)
4166                    return (EFBIG);
4167            error = fget(td, fd, cap_rights_init(&rights, CAP_PWRITE), &fp);
4168            if (error != 0)
4169                    return (error);
4170            if ((fp->f_ops->fo_flags & DFLAG_SEEKABLE) == 0) {
4171                    error = ESPIPE;
4172                    goto out;
4173            }
4174            if ((fp->f_flag & FWRITE) == 0) {
4175                    error = EBADF;
4176                    goto out;
4177            }
4178            if (fp->f_type != DTYPE_VNODE) {
4179                    error = ENODEV;
4180                    goto out;
4181            }
4182            vp = fp->f_vnode;
4183            if (vp->v_type != VREG) {
4184                    error = ENODEV;
4185                    goto out;
4186            }
4187
4188            /* Allocating blocks may take a long time, so iterate. */
4189            for (;;) {
4190                    olen = len;
4191                    ooffset = offset;
4192
4193                    bwillwrite();
4194                    mp = NULL;
```

```
4195                        error = vn_start_write(vp, &mp, V_WAIT | PCATCH);
4196                        if (error != 0)
4197                                break;
4198                        error = vn_lock(vp, LK_EXCLUSIVE);
4199                        if (error != 0) {
4200                                vn_finished_write(mp);
4201                                break;
4202                        }
#ifdef MAC
4204                        error = mac_vnode_check_write(td->td_ucred, fp->f_cred, vp);
4205                        if (error == 0)
#endif
4207                                error = VOP_ALLOCATE(vp, &offset, &len);
4208                        VOP_UNLOCK(vp, 0);
4209                        vn_finished_write(mp);

4211                        if (olen + ooffset != offset + len) {
4212                                panic("offset + len changed from %jx/%jx to %jx/%jx",
4213                                    ooffset, olen, offset, len);
4214                        }
4215                        if (error != 0 || len == 0)
4216                                break;
4217                        KASSERT(olen > len, ("Iteration did not make progress?"));
4218                        maybe_yield();
4219                }
4220  out:
4221        fdrop(fp, td);
4222        return (error);
4223 }

4225 int
4226 sys_posix_fallocate(struct thread *td, struct posix_fallocate_args *uap)
4227 {
4228        int error;

4230        error = kern_posix_fallocate(td, uap->fd, uap->offset, uap->len);
4231        return (kern_posix_error(td, error));
4232 }

4234 /*
4235  * Unlike madvise(2), we do not make a best effort to remember every
4236  * possible caching hint.  Instead, we remember the last setting with
4237  * the exception that we will allow POSIX_FADV_NORMAL to adjust the
4238  * region of any current setting.
4239  */
4240 int
4241 kern_posix_fadvise(struct thread *td, int fd, off_t offset, off_t len,
4242      int advice)
4243 {
4244        struct fadvise_info *fa, *new;
4245        struct file *fp;
4246        struct vnode *vp;
4247        cap_rights_t rights;
4248        off_t end;
4249        int error;

4251        if (offset < 0 || len < 0 || offset > OFF_MAX - len)
4252                return (EINVAL);
4253        switch (advice) {
4254        case POSIX_FADV_SEQUENTIAL:
4255        case POSIX_FADV_RANDOM:
4256        case POSIX_FADV_NOREUSE:
4257                new = malloc(sizeof(*fa), M_FADVISE, M_WAITOK);
4258                break;
4259        case POSIX_FADV_NORMAL:
4260        case POSIX_FADV_WILLNEED:
4261        case POSIX_FADV_DONTNEED:
4262                new = NULL;
4263                break;
4264        default:
4265                return (EINVAL);
4266        }
4267        /* XXX: CAP_POSIX_FADVISE? */
4268        error = fget(td, fd, cap_rights_init(&rights), &fp);
4269        if (error != 0)
4270                goto out;
4271        if ((fp->f_ops->fo_flags & DFLAG_SEEKABLE) == 0) {
4272                error = ESPIPE;
4273                goto out;
4274        }
4275        if (fp->f_type != DTYPE_VNODE) {
4276                error = ENODEV;
4277                goto out;
4278        }
4279        vp = fp->f_vnode;
4280        if (vp->v_type != VREG) {
4281                error = ENODEV;
4282                goto out;
4283        }
4284        if (len == 0)
4285                end = OFF_MAX;
4286        else
4287                end = offset + len - 1;
4288        switch (advice) {
4289        case POSIX_FADV_SEQUENTIAL:
4290        case POSIX_FADV_RANDOM:
4291        case POSIX_FADV_NOREUSE:
4292                /*
```

```
4293                        * Try to merge any existing non-standard region with
4294                        * this new region if possible, otherwise create a new
4295                        * non-standard region for this request.
4296                        */
4297                      mtx_pool_lock(mtxpool_sleep, fp);
4298                      fa = fp->f_advice;
4299                      if (fa != NULL && fa->fa_advice == advice &&
4300                          ((fa->fa_start <= end && fa->fa_end >= offset) ||
4301                          (end != OFF_MAX && fa->fa_start == end + 1) ||
4302                          (fa->fa_end != OFF_MAX && fa->fa_end + 1 == offset))) {
4303                              if (offset < fa->fa_start)
4304                                      fa->fa_start = offset;
4305                              if (end > fa->fa_end)
4306                                      fa->fa_end = end;
4307                      } else {
4308                              new->fa_advice = advice;
4309                              new->fa_start = offset;
4310                              new->fa_end = end;
4311                              fp->f_advice = new;
4312                              new = fa;
4313                      }
4314                      mtx_pool_unlock(mtxpool_sleep, fp);
4315                      break;
4316              case POSIX_FADV_NORMAL:
4317                      /*
4318                       * If a the "normal" region overlaps with an existing
4319                       * non-standard region, trim or remove the
4320                       * non-standard region.
4321                       */
4322                      mtx_pool_lock(mtxpool_sleep, fp);
4323                      fa = fp->f_advice;
4324                      if (fa != NULL) {
4325                              if (offset <= fa->fa_start && end >= fa->fa_end) {
4326                                      new = fa;
4327                                      fp->f_advice = NULL;
4328                              } else if (offset <= fa->fa_start &&
4329                                  end >= fa->fa_start)
4330                                      fa->fa_start = end + 1;
4331                              else if (offset <= fa->fa_end && end >= fa->fa_end)
4332                                      fa->fa_end = offset - 1;
4333                              else if (offset >= fa->fa_start && end <= fa->fa_end) {
4334                                      /*
4335                                       * If the "normal" region is a middle
4336                                       * portion of the existing
4337                                       * non-standard region, just remove
4338                                       * the whole thing rather than picking
4339                                       * one side or the other to
4340                                       * preserve.
4341                                       */
4342                                      new = fa;
4343                                      fp->f_advice = NULL;
4344                              }
4345                      }
4346                      mtx_pool_unlock(mtxpool_sleep, fp);
4347                      break;
4348              case POSIX_FADV_WILLNEED:
4349              case POSIX_FADV_DONTNEED:
4350                      error = VOP_ADVISE(vp, offset, end, advice);
4351                      break;
4352              }
4353  out:
4354          if (fp != NULL)
4355                  fdrop(fp, td);
4356          free(new, M_FADVISE);
4357          return (error);
4358  }
4359
4360  int
4361  sys_posix_fadvise(struct thread *td, struct posix_fadvise_args *uap)
4362  {
4363          int error;
4364
4365          error = kern_posix_fadvise(td, uap->fd, uap->offset, uap->len,
4366              uap->advice);
4367          return (kern_posix_error(td, error));
4368  }
```