# Архитектура операционной системы
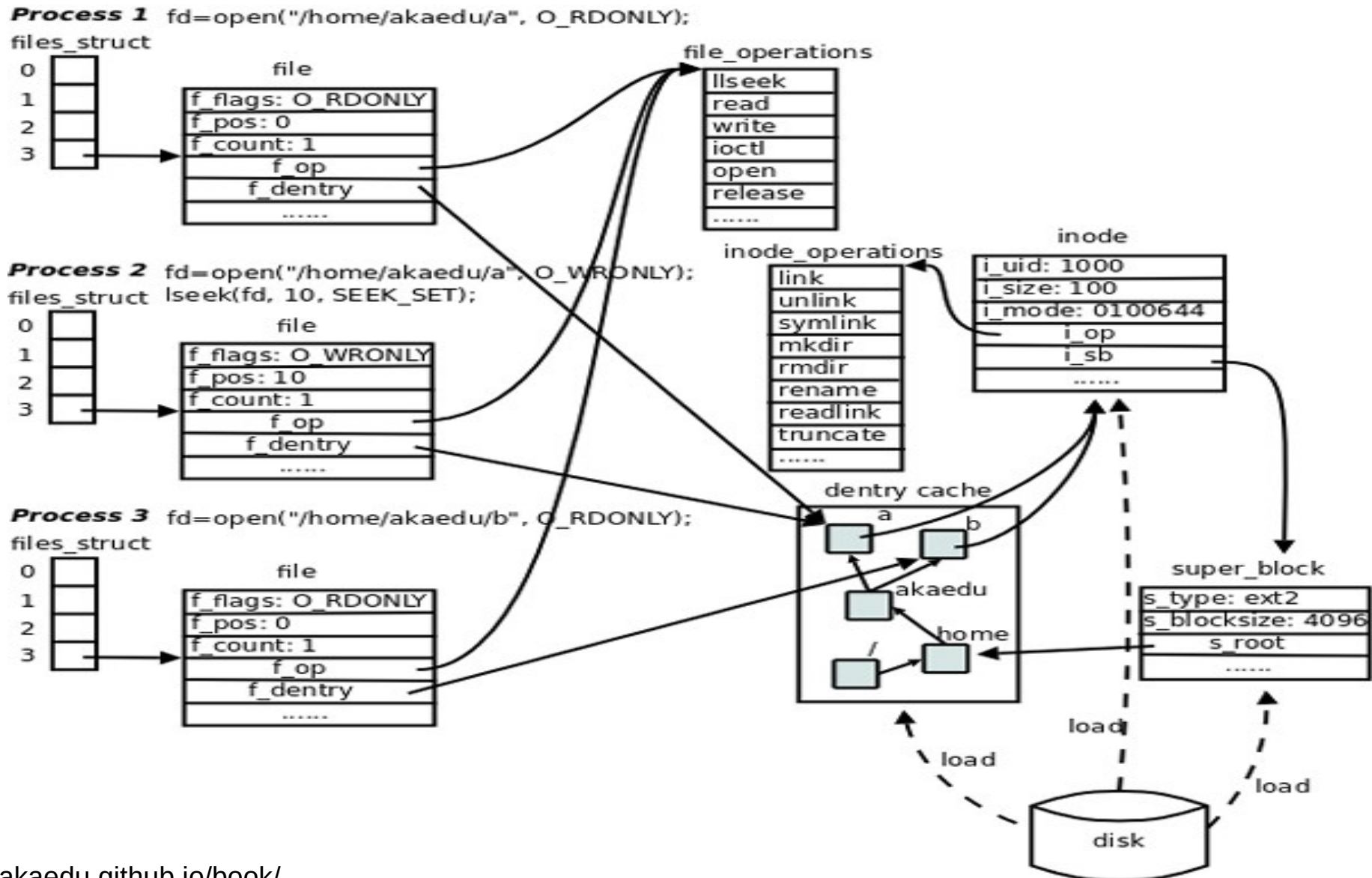**Files, filesystems, vfs**

# Назначение

- Реализация единого интерфейса для построения файловых систем

- Уровень абстракции от конкретной файловой системы

- Стандартизация доступа к ФС на разных носителях

# Структуры данных VFS

- **superblock** – главный (первый) блок файловой системы, описывающий ее параметры (метаинформация о ФС)

- **inode** – индексный узел (метаинформация о файле)

- **dentry** – элемент пути (каталог)

- **file** – объект, представляющий собой открытый файл, связанный с процессом

# Связь структур данных

# file_systems →* file_system_type

```c
struct file_system_type {
        const char *name;
        int fs_flags;
#define FS_REQUIRES_DEV         1
#define FS_BINARY_MOUNTDATA     2
#define FS_HAS_SUBTYPE          4
#define FS_USERNS_MOUNT         8          /* Can be mounted by userns root */
#define FS_RENAME_DOES_D_MOVE   32768   /* FS will handle d_move() during rename() internally. */
        struct dentry *(*mount) (struct file_system_type *, int,
                        const char *, void *);
        void (*kill_sb) (struct super_block *);
        struct module *owner;
        struct file_system_type * next;
        struct hlist_head fs_supers;

        struct lock_class_key s_lock_key;
        struct lock_class_key s_umount_key;
        struct lock_class_key s_vfs_rename_key;
        struct lock_class_key s_writers_key[SB_FREEZE_LEVELS];

        struct lock_class_key i_lock_key;
        struct lock_class_key i_mutex_key;
        struct lock_class_key i_mutex_dir_key;
};
```

# Структуры данных VFS

- **superblock** – главный (первый) блок файловой системы, описывающий ее параметры (метаинформация о ФС)

- **inode** – индексный узел (метаинформация о файле)

- **dentry** – элемент пути (каталог)

- **file** – объект, представляющий собой открытый файл, связанный с процессом

# current → namespace → * vfsmount

```c
struct vfsmount {
        struct dentry *mnt_root;        /* root of the mounted tree */
        struct super_block *mnt_sb;     /* pointer to superblock */
        int mnt_flags;
} __randomize_layout;

struct file; /* forward dec */
```

```c
struct super_block {
        struct list_head        s_list;          /* Keep this first */
        dev_t                   s_dev;           /* search index; _not_ kdev
        unsigned char           s_blocksize_bits;
        unsigned long           s_blocksize;
        loff_t                  s_maxbytes;      /* Max file size */
        struct file_system_type *s_type;
        const struct super_operations   *s_op;
        const struct dquot_operations   *dq_op;
        const struct quotactl_ops       *s_qcop;
        const struct export_operations *s_export_op;
        unsigned long           s_flags;
        unsigned long           s_iflags;        /* internal SB_I_* flags */
        unsigned long           s_magic;
        struct dentry           *s_root;
        struct rw_semaphore     s_umount;
        int                     s_count;
        atomic_t                s_active;
#ifdef CONFIG_SECURITY
        void                    *s_security;
#endif
        const struct xattr_handler **s_xattr;
```

```c
1777    struct super_operations {
1778            struct inode *(*alloc_inode)(struct super_block *sb);
1779            void (*destroy_inode)(struct inode *);
1780
1781            void (*dirty_inode) (struct inode *, int flags);
1782            int (*write_inode) (struct inode *, struct writeback_control *wbc);
1783            int (*drop_inode) (struct inode *);
1784            void (*evict_inode) (struct inode *);
1785            void (*put_super) (struct super_block *);
1786            int (*sync_fs)(struct super_block *sb, int wait);
1787            int (*freeze_super) (struct super_block *);
1788            int (*freeze_fs) (struct super_block *);
1789            int (*thaw_super) (struct super_block *);
1790            int (*unfreeze_fs) (struct super_block *);
1791            int (*statfs) (struct dentry *, struct kstatfs *);
1792            int (*remount_fs) (struct super_block *, int *, char *);
1793            void (*umount_begin) (struct super_block *);
1794
1795            int (*show_options)(struct seq_file *, struct dentry *);
1796            int (*show_devname)(struct seq_file *, struct dentry *);
1797            int (*show_path)(struct seq_file *, struct dentry *);
1798            int (*show_stats)(struct seq_file *, struct dentry *);
1799    #ifdef CONFIG_QUOTA
```

```
191
192    struct export_operations {
193            int (*encode_fh)(struct inode *inode, __u32 *fh, int *max_len,
194                             struct inode *parent);
195            struct dentry * (*fh_to_dentry)(struct super_block *sb, struct fid *fid,
196                             int fh_len, int fh_type);
197            struct dentry * (*fh_to_parent)(struct super_block *sb, struct fid *fid,
198                             int fh_len, int fh_type);
199            int (*get_name)(struct dentry *parent, char *name,
200                             struct dentry *child);
201            struct dentry * (*get_parent)(struct dentry *child);
202            int (*commit_metadata)(struct inode *inode);
203
204            int (*get_uuid)(struct super_block *sb, u8 *buf, u32 *len, u64 *offset);
205            int (*map_blocks)(struct inode *inode, loff_t offset,
206                             u64 len, struct iomap *iomap,
207                             bool write, u32 *device_generation);
208            int (*commit_blocks)(struct inode *inode, struct iomap *iomaps,
209                             int nr_iomaps, struct iattr *iattr);
210    };
211
212    extern int exportfs_encode_inode_fh(struct inode *inode, struct fid *fid,
213                             int *max_len, struct inode *parent);
214    extern int exportfs_encode_fh(struct dentry *dentry, struct fid *fid,
```

```
847
848    struct file {
849            union {
850                    struct llist_node          fu_llist;
851                    struct rcu_head            fu_rcuhead;
852            } f_u;
853            struct path                  f_path;
854            struct inode                 *f_inode;      /* cached value */
855            const struct file_operations  *f_op;
856
857            /*
858             * Protects f_ep_links, f_flags.
859             * Must not be taken from IRQ context.
860             */
861            spinlock_t                   f_lock;
862            enum rw_hint                 f_write_hint;
863            atomic_long_t                f_count;
864            unsigned int                 f_flags;
865            fmode_t                      f_mode;
866            struct mutex                 f_pos_lock;
867            loff_t                       f_pos;
868            struct fown_struct           f_owner;
869            const struct cred            *f_cred;
870            struct file_ra_state         f_ra;
```

```
561     /*
562      * Keep mostly read-only and often accessed (especially for
563      * the RCU path lookup and 'stat' data) fields at the beginning
564      * of the 'struct inode'
565      */
566     struct inode {
567             umode_t                 i_mode;
568             unsigned short          i_opflags;
569             kuid_t                  i_uid;
570             kgid_t                  i_gid;
571             unsigned int            i_flags;
572
573     #ifdef CONFIG_FS_POSIX_ACL
574             struct posix_acl        *i_acl;
575             struct posix_acl        *i_default_acl;
576     #endif
577
578             const struct inode_operations   *i_op;
579             struct super_block      *i_sb;
580             struct address_space    *i_mapping;
581
582     #ifdef CONFIG_SECURITY
583             void                    *i_security;
584     #endif
645             const struct file_operations    *i_fop; /* former ->i_op->default_file_
```

```c
1664   struct file_operations {
1665           struct module *owner;
1666           loff_t (*llseek) (struct file *, loff_t, int);
1667           ssize_t (*read) (struct file *, char __user *, size_t, loff_t *);
1668           ssize_t (*write) (struct file *, const char __user *, size_t, loff_t *)
1669           ssize_t (*read_iter) (struct kiocb *, struct iov_iter *);
1670           ssize_t (*write_iter) (struct kiocb *, struct iov_iter *);
1671           int (*iterate) (struct file *, struct dir_context *);
1672           int (*iterate_shared) (struct file *, struct dir_context *);
1673           unsigned int (*poll) (struct file *, struct poll_table_struct *);
1674           long (*unlocked_ioctl) (struct file *, unsigned int, unsigned long);
1675           long (*compat_ioctl) (struct file *, unsigned int, unsigned long);
1676           int (*mmap) (struct file *, struct vm_area_struct *);
1677           int (*open) (struct inode *, struct file *);
1678           int (*flush) (struct file *, fl_owner_t id);
1679           int (*release) (struct inode *, struct file *);
1680           int (*fsync) (struct file *, loff_t, loff_t, int datasync);
1681           int (*fasync) (int, struct file *, int);
1682           int (*lock) (struct file *, int, struct file_lock *);
1683           ssize_t (*sendpage) (struct file *, struct page *, int, size_t, loff_t
1684           unsigned long (*get_unmapped_area)(struct file *, unsigned long, unsigne
1685           int (*check_flags)(int);
1686           int (*flock) (struct file *, int, struct file_lock *);
```

```
1703
1704    struct inode_operations {
1705            struct dentry * (*lookup) (struct inode *,struct dentry *, unsigned int
1706            const char * (*get_link) (struct dentry *, struct inode *, struct delay
1707            int (*permission) (struct inode *, int);
1708            struct posix_acl * (*get_acl)(struct inode *, int);
1709
1710            int (*readlink) (struct dentry *, char __user *,int);
1711
1712            int (*create) (struct inode *,struct dentry *, umode_t, bool);
1713            int (*link) (struct dentry *,struct inode *,struct dentry *);
1714            int (*unlink) (struct inode *,struct dentry *);
1715            int (*symlink) (struct inode *,struct dentry *,const char *);
1716            int (*mkdir) (struct inode *,struct dentry *,umode_t);
1717            int (*rmdir) (struct inode *,struct dentry *);
1718            int (*mknod) (struct inode *,struct dentry *,umode_t,dev_t);
1719            int (*rename) (struct inode *, struct dentry *,
1720                            struct inode *, struct dentry *, unsigned int);
1721            int (*setattr) (struct dentry *, struct iattr *);
1722            int (*getattr) (const struct path *, struct kstat *, u32, unsigned int)
1723            ssize_t (*listxattr) (struct dentry *, char *, size_t);
1724            int (*fiemap)(struct inode *, struct fiemap_extent_info *, u64 start,
1725                            u64 len);
1726            int (*update_time)(struct inode *, struct timespec *, int);
```

```
89    struct dentry {
90            /* RCU lookup touched fields */
91            unsigned int d_flags;           /* protected by d_lock */
92            seqcount_t d_seq;               /* per dentry seqlock */
93            struct hlist_bl_node d_hash;    /* lookup hash list */
94            struct dentry *d_parent;        /* parent directory */
95            struct qstr d_name;
96            struct inode *d_inode;          /* Where the name belongs to - NULL is
97                                             * negative */
98            unsigned char d_iname[DNAME_INLINE_LEN];    /* small names */
99
100           /* Ref lookup also touches following */
101           struct lockref d_lockref;       /* per-dentry lock and refcount */
102           const struct dentry_operations *d_op;
103           struct super_block *d_sb;       /* The root of the dentry tree */
104           unsigned long d_time;           /* used by d_revalidate */
105           void *d_fsdata;                 /* fs-specific data */
106
107           union {
108                   struct list_head d_lru;         /* LRU list */
109                   wait_queue_head_t *d_wait;      /* in-lookup ones only */
110           };
111           struct list_head d_child;       /* child of parent list */
```

```c
 67          *          unregistered.
 68          */
 69
 70     int register_filesystem(struct file_system_type * fs)
 71     {
 72                 int res = 0;
 73                 struct file_system_type ** p;
 74
 75                 BUG_ON(strchr(fs->name, '.'));
 76                 if (fs->next)
 77                         return -EBUSY;
 78                 write_lock(&file_systems_lock);
 79                 p = find_filesystem(fs->name, strlen(fs->name));
 80                 if (*p)
 81                         res = -EBUSY;
 82                 else
 83                         *p = fs;
 84                 write_unlock(&file_systems_lock);
 85                 return res;
 86     }
 87
 88     EXPORT_SYMBOL(register_filesystem);
 89
 90     /**
```