

# Distributed file systems

Johan Montelius

KTH

HT15

Functionality:

- persistent storage of files, create and delete
- manipulating a file, read and write operations
- authorization, who is allowed to do what
- a naming service

*The mapping of names to files is quite separate from the rest of the system.*

1 / 31

2 / 31

directory	map from name to identifier
file module	locate file, harder in distributed systems
access control	interacts with authentication system
file operations	read and write operations
block operations	
device operations	

3 / 31

- a sequence of bytes
- attributes, associated meta-data
  - size and type
  - owner and permissions
  - author
  - created, last written
  - icons, fonts, presentation...

4 / 31

- `create(name, mode)` returns a file-descriptor
- `open(name, mode)` returns a file-descriptor
- `close(fd)`
  
- `unlink(name)`
- `link(name, name)`

Can we separate the name service from the file operations?

Operating system operations are not always directly available from a high level language.

Buffering of write operations to reduce the number of system calls.

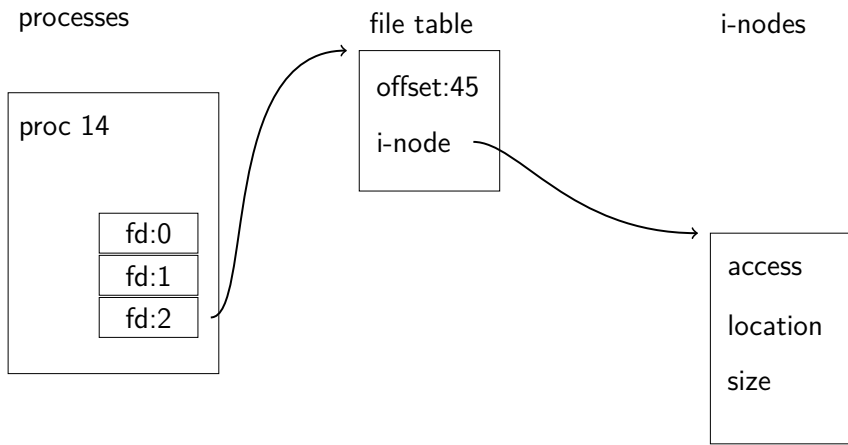
- `read(fd, buffer, n)` returns the number of bytes actually read
- `write(fd, buffer, n)` returns the number of bytes actually written
- `lseek(fd, offset, set/cur/end)` sets the current position in the file
  
- `stat(name, buffer)` reads the file attributes

A process holds a set of open *file descriptors*, each descriptor holds a pointer to a table of open files.

The file table entries holds a *position* and a pointer to an *inode* (index node).

The inode holds information about where file blocks are allocated.

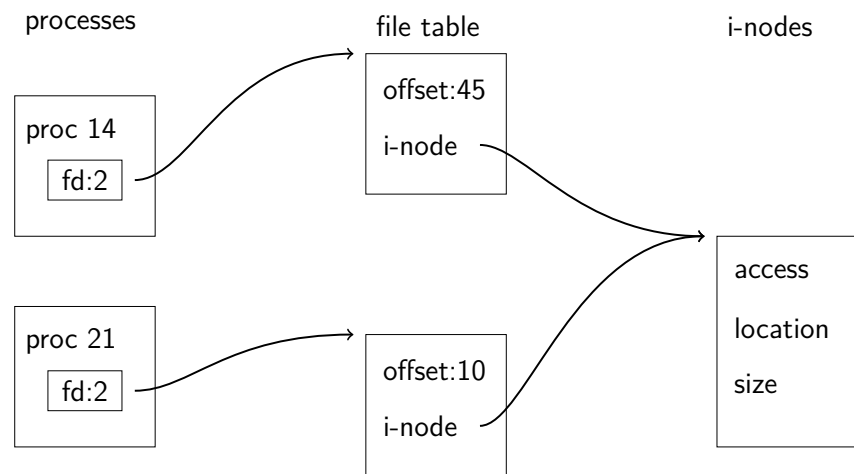
## descriptors, table entries and i-nodes



All *read*, *write* and *lseek* operations will change the current position in the table entry.

9 / 31

## two processes open the same file



Two processes that open the same file will have independent file table entries.

10 / 31

## Nota bene

All threads in a process (or even if process is *forked*) share the same file descriptors and thus share the file table entry.

11 / 31

## one-copy semantics

Most file systems give us a one-copy semantics

- we expect operations to be visible by everyone and that everyone see the same file
- if I tell you that the file has been modified the modification should be visible

12 / 31

## an architecture of a distributed file system

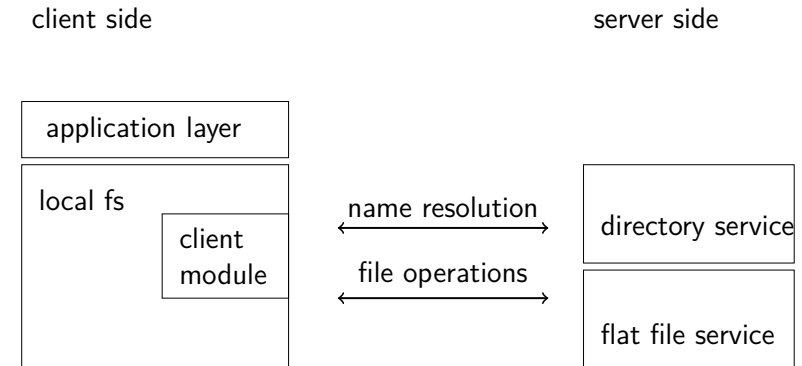
Let's define the requirements we have on a distributed file system.

- transparency - no difference between local and remote files
  - access: same set operations
  - location: same name space
  - mobility: allowed to move files without changing client side
  - performance: close to a non-distributed system
- concurrency - simultaneous operations by several clients
- heterogeneity - not locked in to a particular operating system
- fault tolerance - independent of clients, restartable etc
- consistency - one-copy semantics... or?
- security - access control, who is allowed to do what

13 / 31

## distributed architecture

Separate the directory service from the file service.



14 / 31

## the directory service

The directory service - what operations do we need?

- lookup a *file identifier* given a *name* and *directory*
- link a *name* to a *file identifier* in a *directory*
- remove a link
- list all *names* in a *directory*

*The directory service does not create nor manipulate files.*

15 / 31

## the file service

What operations should be provided?

- create a file and allocate a *file identifier*
- delete a file
- read from a file identified by a *file identifier*
- write a number of bytes to a file identified by a *file identifier*

Do we need a open operation?

- What does open do in Unix?
- What do we need if we don't have an open operation?
- What would the benefit be?

16 / 31

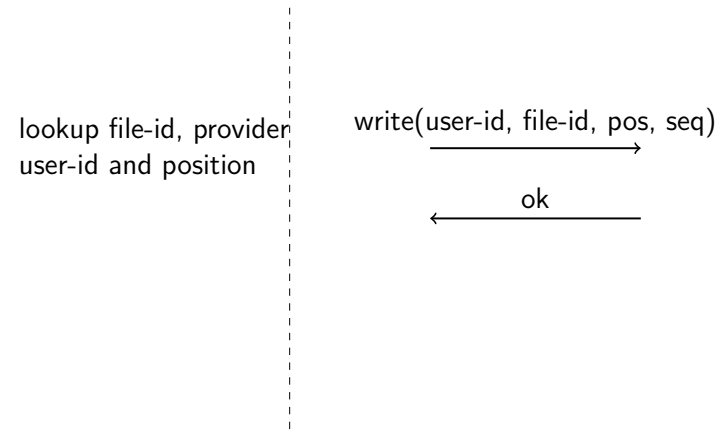
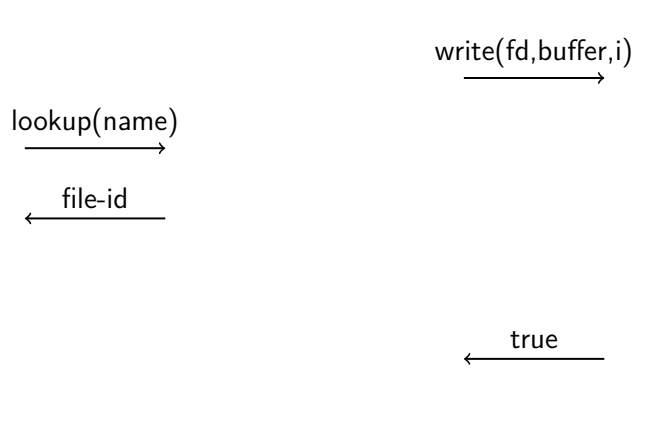
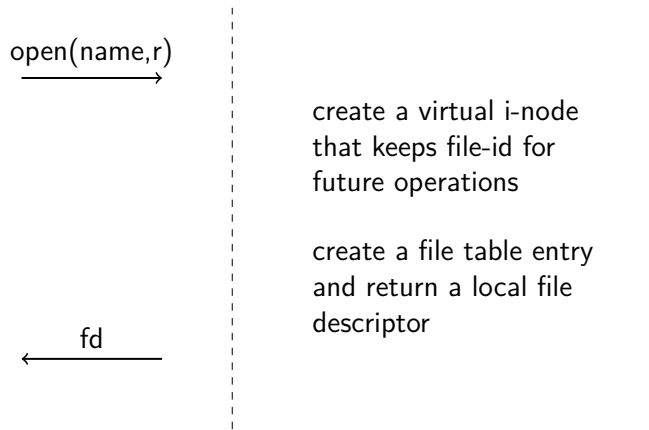
- What are the benefits of a *stateless server*?
- What are the benefits of a *stateless server*?
- How can we maintain a *session state* while keeping the server stateless?

In Unix, permissions are checked **when a file is opened** and access to the file can then be done without security control.

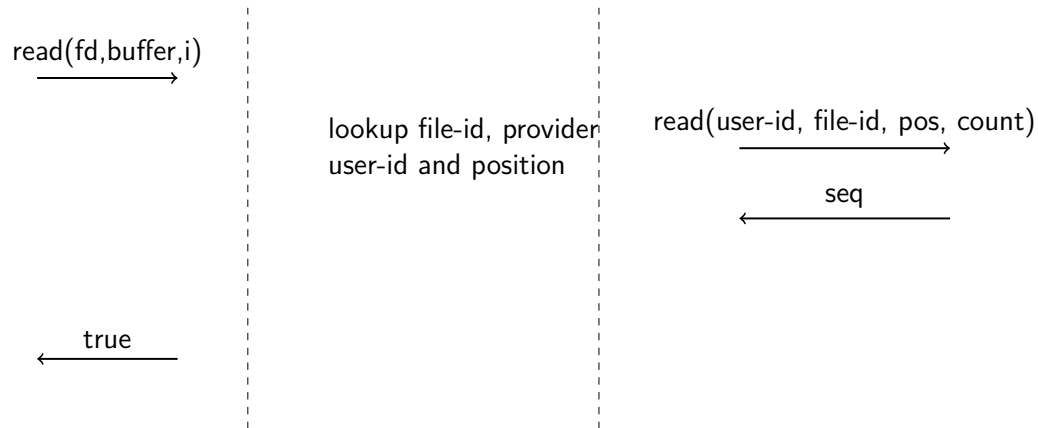
If we do not have an *open operation*, how can we perform authentication and authorization control?

Client interface - open

Client interface - write

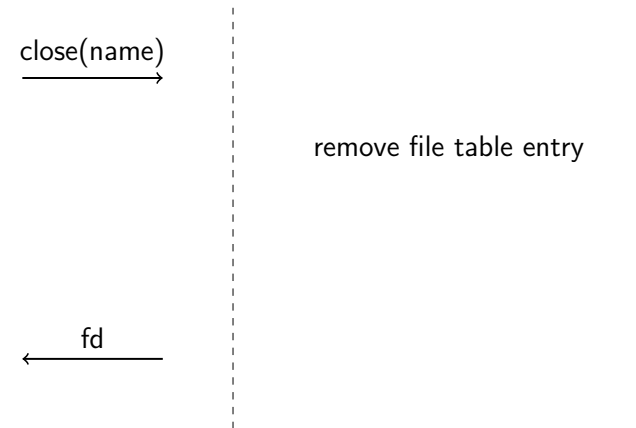


## Client interface - read



21 / 31

## Client interface - close



22 / 31

## Performance

Everything would be fine, if it was not for performance.

Keep a local copy of the file at the client side.

23 / 31

## Caching - options

Reading from a file: how do we know it is the most current?

- check validity when reading
- ... if you haven't done so in a while
- server should invalidate copy

Writing to a file:

- write-through: write to cache and server
- write-back : write to cache only
- write-around: write only to server

*caching could break the one-copy semantics*

24 / 31

- developed by Sun, 1984
- targeting department networks
- implemented using RPC (Open Network Computing)
- public API: RFC 1094, 1813, 3530
- originally used UDP, later versions have support for TCP to improve performance over WAN
- mostly used with UNIX systems but client on all platforms available

### Reading from a file:

- first read will copy a segment (8K bytes) to the client
- the copy *valid* of a time  $t$  (3-30 seconds)
- if more time has elapsed, the validity is checked again

*the server is stateless*

### Writing to a file:

- write-back : write to the cache only
- schedule written segment to be copied to server
- segment copied on timeout or when file is closed (sync)

- developed by Carnegie Mellon University
- clients for most platforms, OpenAFS (from IBM), Arla (a KTH implementation)
- used mainly in WAN (Internet) where the overhead of NFS would be prohibitive
- caching of whole files and infrequent sharing of writable files

### Reading from a file:

- copy the whole file from server (or 64kbyte)
- receive a *call-back promise*
- file is valid if promise exists and is not too old (minutes)

### Writing to a file:

- write-back : write to the cache only
- file copied to server when closed (sync)
- server will invalidate all existing promises

- Service Message Block (SMB) was originally developed by IBM but then modified by Microsoft, now also under the name Common Internet File System (CIFS).
- not only file sharing but also name servers, printer sharing etc.
- Samba is an open source reimplementation of SMB by Andrew Tridgell

- SMB uses *client locks* to solve cache consistency
- a client can open a file and lock it; all read and write operations in client cache
- a read only lock will allow multiple clients to cache and read a file
- Locks can be revoked by the server forcing the client to flush any changes
- in a unreliable or high latency network, locking can be dangerous and counter productive

29 / 31

30 / 31

## Summary

- separate directory service from file service
- maintain a view of only one file, one-copy semantics
- caching is key to performance but could make the one-copy view hard to maintain

31 / 31