

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



- Understanding the need for Distributed File Systems (DFS)

Revisiting the basics of Unix File Systems

- Understanding the key requirements for DFS
- Exp
- Cas

- Reading: Distributed Systems: Concepts and Design by George Coulouris (5th edition). Chapter 12. Se

Assignment Project Exam Help

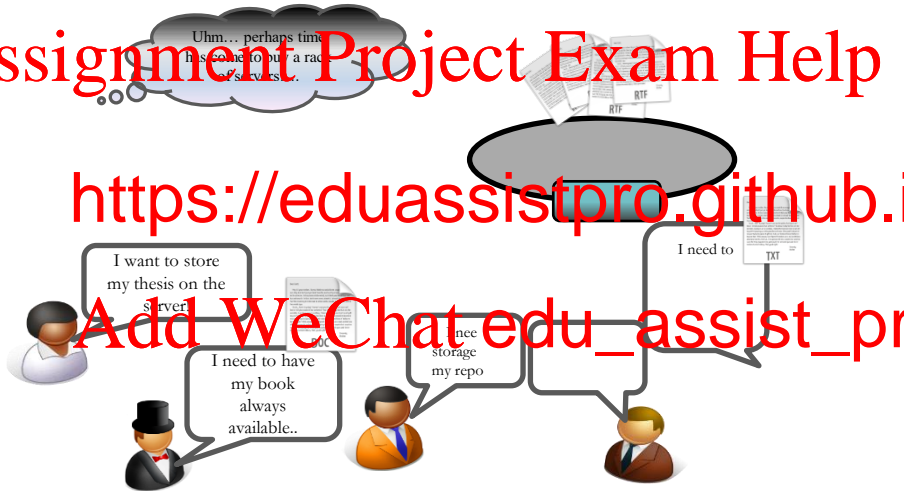
<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



Assignment Project Exam Help

<https://eduassistpro.github.io>

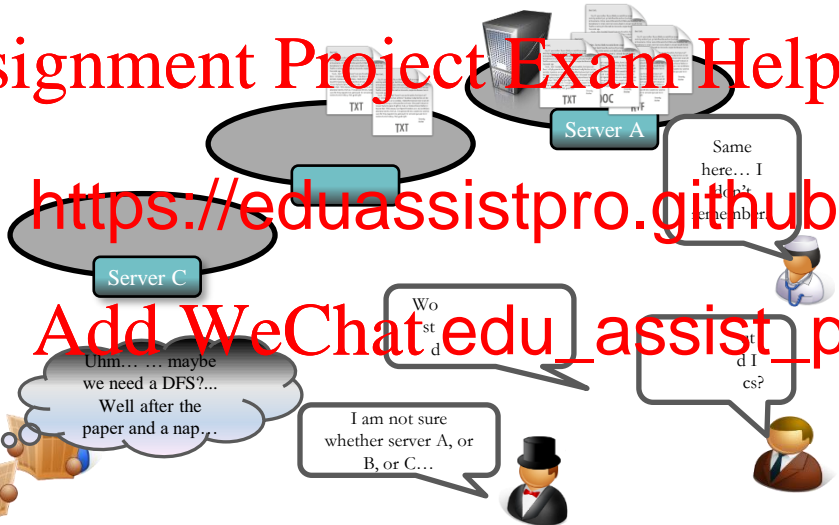




Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro





Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



It is reliable, fault tolerant, highly available, location transparent.... I hope I can finish my newspaper now...

Nice... my boss will promote me!

have to remember which server I stored the data into...

I can access folders from anywhere..



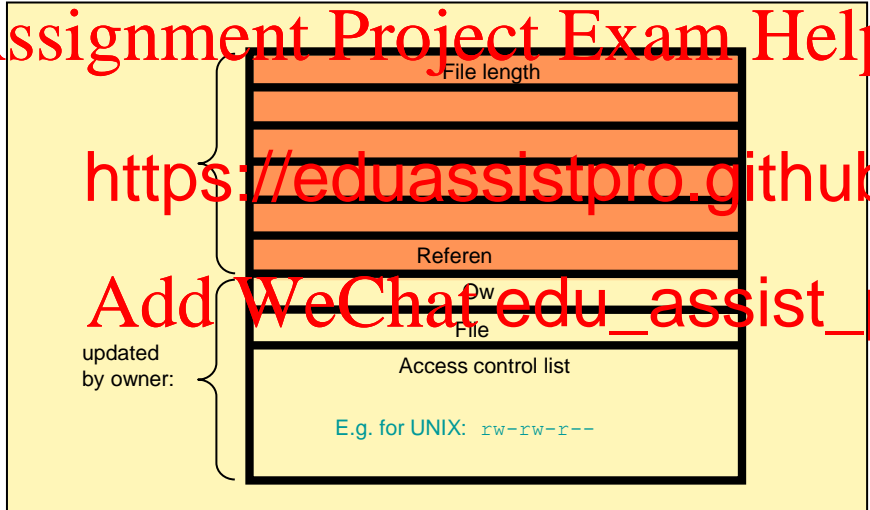


	<i>Sharing</i>	<i>Persis- tence</i>	<i>Distributed cache replicas</i>	<i>Consistency maintenance</i>	<i>Example</i>
Main memory	×	×	×	1	RAM
File sys	×	✓	×		file system
Distribu	✓	✓	✓		FS
Web	✓	✓	✓	×	server
Distributed shared memory	✓	×	✓	✓	Ivy (Ch. 16)
Remote objects (RMI/GRB)	✓	×	×	×	
Persistent object store	✓	✓	×	×	ersistent e
Peer-to-peer storage store	✓	✓	✓	2	OceanStore

Types of consistency between copies: 1 - strict one-copy consistency ✓ - approximate/slightly weaker guarantees
X - no automatic consistency 2 – considerably weaker guarantees



- Assignment Project Exam Help**
- Files contain both *data* and *attributes*. Data is usually in the form of a sequence of bytes and can be accessed and modified. Attributes include things like file name, file size, file type, file permissions, and access control. Files have a unique identity and are used to organize data.
 - Files have a hierarchical naming scheme, or the file path, which is a concatenation of the directory names and the file name. Files can be organized into a hierarchy of other files.
- https://eduassistpro.github.io**
- Add WeChat edu_assist_pro**

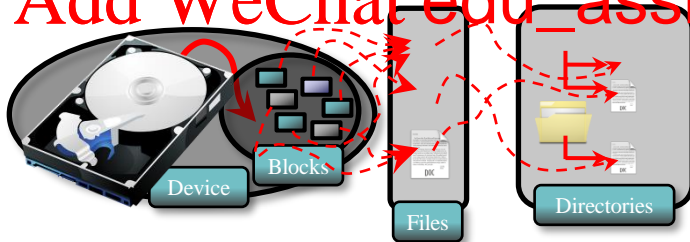




Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr





filedes = *open(name, mode)*
filedes = *creat(name, mode)*

Open an existing file with the given *name*.
Creates a new file with the given *name*.

Both operations deliver a file descriptor referencing the open file. The *mode* is *read*, *write* or both.

status

count =

count =

count bytes transferred
from buffer.
to buffer.

and advance the read-write pointer.

pos = *lseek(filedes, offset, whence)*

Moves the read-write
depending on *whence*

status = *unlink(name)*

Removes the file *name* if file
has no other names, it

status = *link(name1, name2)*

Adds a new name (*name2*) for a file (*name1*).

status = *stat(name, buffer)*

Gets the file attributes for file *name* into *buffer*.



Write a simple C program to copy a file using the UNIX file system operations.

```
#define BUFSIZE 1024
void copyfile(char* oldfile, char* newfile)
{
    char buf[BUFSIZE]; int i,n=1, fdold, fdnew;

    fdold = open(oldfile, O_RDONLY);
    if (fdold < 0) {
        printf("Copyfile: couldn't open %s\n", oldfile);
        return;
    }
    fdnew = open(newfile, O_WRONLY | O_CREAT | O_TRUNC, 0666);
    if (fdnew < 0) {
        printf("Copyfile: couldn't open %s\n", newfile);
        return;
    }
    while (n = read(fdold, buf, BUFSIZE)) {
        if (write(fdnew, buf, n) != n) {
            printf("Copyfile: write error\n");
            return;
        }
    }
    close(fdold);
    close(fdnew);
}

main(int argc, char **argv) {
    copyfile(argv[1], argv[2]);
}
```

PTENR)

```
if(write(
}
close(fdold);
```

```
else printf("Copyfile: couldn't open f
```

```
main(int argc, char **argv) {
    copyfile(argv[1], argv[2]);
}
```

- A *file system* provides a convenient programming interface for disk storage along with features such as access control and file-locking that allows file sharing.
- A basic *distributed file system* emulates the same functionality as a (non-distributed) remote computer.
- A *file* exactly as they do local ones, allowing users to access their files from any computer in an intranet.
- Hosts that provide a file service can be a wide range of other services in an organization, e.g., for the web services and email services. This further facilitates management of the persistent storage, including backups and archiving.

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Transparency :

Assignment Project Exam Help

Access transparency – Client programs should be unaware of the distribution of files. Same API is used for accessing local and remote files and so programs written to operate on local

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

- *Mobility transparency* – Client services do not need to change one place to another.
- *Performance transparency* – perform satisfactorily while the specified range.
- *Scaling transparency* – The service can be expanded by incremental growth to deal with a wide range of loads and network sizes.

a uniform
ent
here the

istration
from

tinue to
within a



Concurrent file updates : Multiple clients' updates to files should not interfere with each other. Policies should be manageable.

File replication : Each file can have multiple copies distributed over

processing

Hardware

<https://eduassistpro.github.io>

should not

require the client or server to
operating system dependencies

Fault tolerance : Transient communication failure or file

corruption. Servers can use at
semantics or the simpler at-least-once semantics with
idempotent operations. Servers can also be *stateless*.



Assignment Project Exam Help

Consistency: Multiple, possibly concurrent, access to a file should see a consistent representation of that file, i.e. differences in the the file should be

<https://eduassistpro.github.io>

Security: Client requests should be authenticated and encrypted.

Efficiency: Should be of a comparable level to conventional file systems.

Add WeChat edu_assist_pro



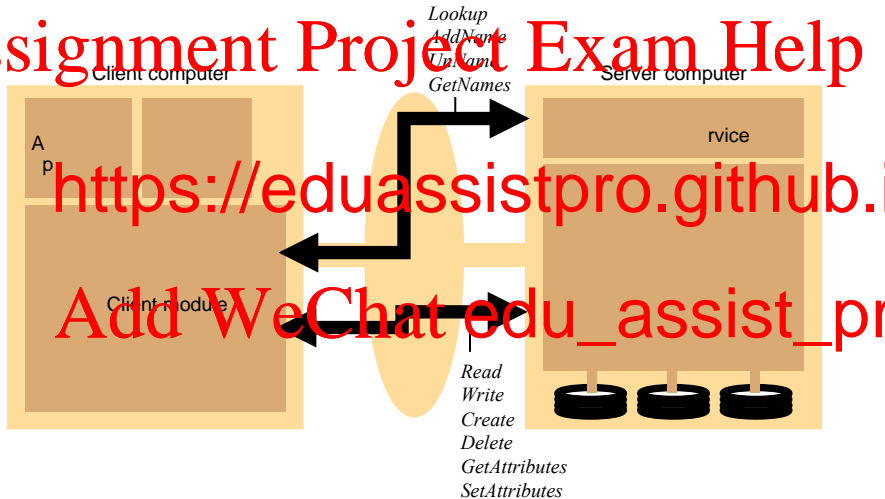
- **Flat file service:** The flat file service is concerned with implementing operations on the contents of files. A *unique file identifier* (UFID) is given to the flat file service to refer to the file to be operated on. The UFID is unique over all the files in the distributed system. The flat file service creates a new UFID for each new file that it
- **Directory service:** The directory service is concerned with implementing operations on the contents of files. A *unique file identifier* (UFID) is given to the directory service to refer to the file to be operated on. The UFID is unique over all the files in the distributed system. The directory service creates a new UFID for each new file that it
- **Client module:** The client module integrates the flat file service and the directory service to provide whatever application programs the application programs. The client module servers. It can also cache data in order to improve performance.

text names
and delete
the flat file

file
d by
file



Assignment Project Exam Help





Assignment Project Exam Help

Read(UFID, i, n) \rightarrow Data

Reads up to n items from position i in the file.

Write(U

the file is

Create()

Delete()

GetAttributes(UFID) \rightarrow Attr

Returns the file attributes for the file.

SetAttributes(UFID,Attr)

Sets the file attribute

<https://eduassistpro.github.io>
Add WeChat edu_assist_pr



Difference with UNIX interface:

- Recall that the UNIX interface shown earlier requires that the UNIX file system maintains state, i.e., a file pointer, that is manipulated during reads
- The flat file service is stateless – the file service does not maintain state and can be restarted after a failure and resume operations for clients or the server to restore any state.
- Also note that UNIX files require an explicit open command before they can be accessed, while files in the flat file service can be accessed immediately.

Assignment Project Exam Help
<https://eduassistpro.github.io>
Add WeChat: edu_assist_pro



Assignment Project Exam Help

- The service needs to authenticate the RPC caller and needs to ensure that illegal operations are not performed, e.g., that UFIDs are legal and that files are
- The service needs to ensure that the client's identity is not broken

Two ways to do this:

1. An access check can be made whenever a file is requested, and the results can be encoded in the form of a client for submission to the flat file server.
2. A user identity can be submitted with every client request, and access checks can be performed by the flat file server for every file operation.



- The primary purpose of the directory service is to provide a translation from file names to UIDs. An abstract directory service interface is shown on the next slide.
- The directory service provides mappings between flat file server and the flat file server.
- A hierarchical file system can be built up from the root directory has name "/" and the contents "usr", "home", "etc", which themselves are directories. For example, the files in the "usr" directory proceed through the path to the file or directory at the end.



Assignment Project Exam Help

Lookup(Dir, Pattern) → Name Returns all the patterns in the given directory.

AddName(Dir, Name) Adds the file name to the directory.

UnName(Dir, Name) Remove the file name from the directory.

GetNames(Dir, Pattern) → Name Returns all the patterns in the directory with the pattern.

Add WeChat edu_assist_pro



Assignment Project Exam Help

A *file group* is a collection of files located on a given server. A server may hold several file groups and file groups can be moved between servers, but a

- File group files are identified by a common file ID (UFID) by which the file was created (16 bits). This allows the file to be moved to a different server without conflicting with files on other servers. For example, if a file is moved from one server to another, the file's UFID remains the same. The file service needs to maintain a mapping of UFIDs to servers. This can be cached at the client module.

file group id:

32 bits

16 bits

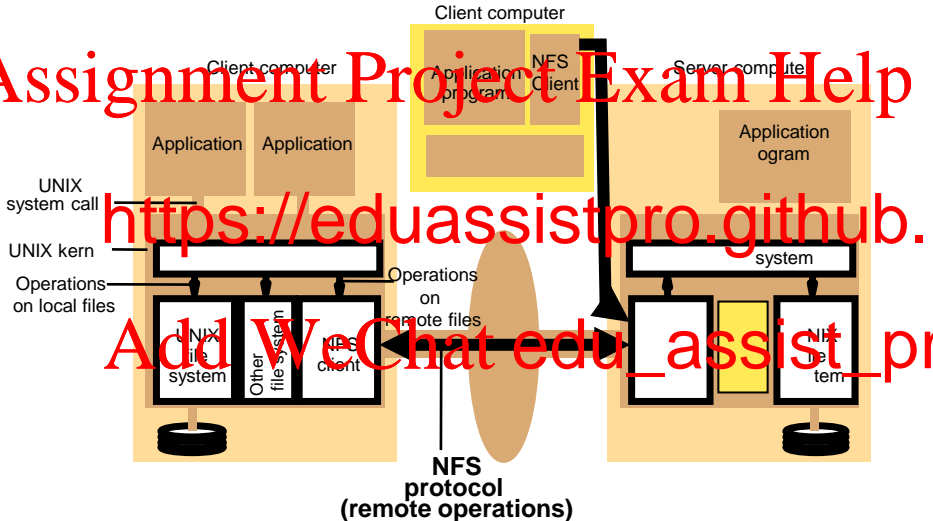
IP address	date
------------	------



Assignment Project Exam Help

- The Sun Network File System (NFS) follows the abstract system shown earlier
- There are two main components: the NFS server and the NFS client. The NFS client performs operations on the remote file store.
- We consider a UNIX implementation.
- The NFS client makes requests to the NFS server.

<https://eduassistpro.github.io>
Add WeChat: edu_assist_pro





- UNIX uses a *virtual file system* (VFS) to provide transparent access to any number of different file systems. The VFS is integrated in the same way.

The VFS maintains a VFS structure for each filesystem in use. The VFS struct .e., it

comb

- The V indicator as to whether the file is local or

- If the file is local, then the v-node contains the UNIX file system.

- If the file is remote then the v-node contains the i-node

file handle which is a combination of *filesystem identifier*, i-node number and whatever else the NFS server needs to identify the file.



Assignment Project Exam Help

The NFS client is integrated within the kernel so that:

- user programs can access files via UNIX system calls without recompilation or reloadi
- a singl
- the en
- retain

ared cache;

er can be
s.

The client transfers blocks of files from the
and caches them, sharing the same buffer c
input-output system. Since several hosts ma
remote file, caching presents a problem of c

st

<https://eduassistpro.github.io>
Add WeChat: edu_assist_pr



- The NFS server interface integrates both the directory and file operations in a single service. The creation and insertion of file names in directories is performed by a single create operation, which takes the text name of the new file and file handle for the target directory as arguments.

- The p system
primi

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



- Each server maintains a file that describes which parts of the local filesystems that are available for remote mounting

```
tawfiq@http:~$ cat /etc/exports
```

```
# /et
```

ay be

```
# exp
```

```
/ stor
```

- In the above example all hosts on the su filesystem directory /store with read an

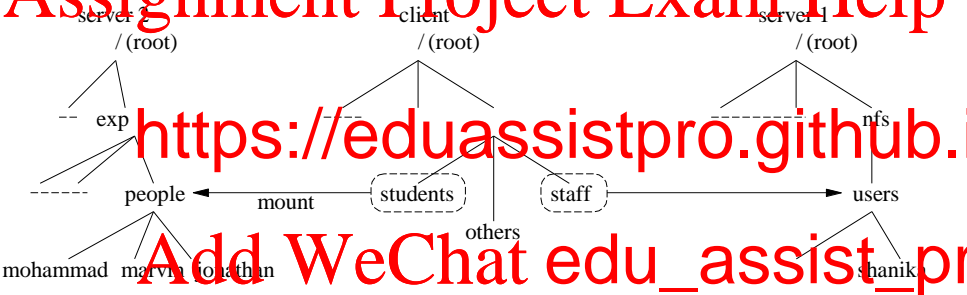
- A *hard-mounted* filesystem will block on

ess IS

complete. A *soft-mounted* filesystem will retry a few times and then return an error to the calling process.



Assignment Project Exam Help



Example NFS mounting in two different file systems



In conventional UNIX systems:

data read from the disk or pages are retained in a main memory buffer cache and are evicted when the buffer space is required for other pages. Accesses to cached data does not require a disk access.

- *Read-ahead*

that h

- *Delayed Write*

when they have been both modified and evicted. A UNIX sync operation flushes modified pages to disk every 30 seconds. This is not the case in a distributed system, on a single host, because there is only one cache. In a distributed system, the client must bypass the cache.



- Use of the cache at the server for client reads does not introduce any problems. However, use of the cache for writes requires special care to ensure that client can be confident that the writes are persistent, especially in the event of a server crash.

There are

- **Write-through** – data is written to cache and is written to disk when a commit operation is received for the data. A reply to the commit operation is written to disk. is increases
- The first option is poor when the server receives requests for the same data. It however save
- The second option uses more network bandwidth and may lead to uncommitted data being lost. However, it receives the full benefit of the cache.

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



- The NFS client also caches data reads, writes, attributes and directory operations in order to reduce network I/O.

Caching at the client introduces the cache consistency problem since now there is a cache at the client and the server, and there may be more than one c

- Note another write on the second read operation being incorrect.

- In NFS, clients poll the server to check fo

Assignment Project Exam Help
<https://eduassistpro.github.io>
Add WeChat edu_assist_pr



- Let T_c be the time when a cache block was last validated by the client. Let T_m be the time when a block was last modified.
- A cache block is said to be valid at time T if (i) $T - T_c < t$ where t is a given *freshness interval*, or (ii) the value of T_m at the client matches the value at the server

<https://eduassistpro.github.io>

- A small value for t leads to a close approximation of one-copy consistency, at the cost of greater network I/O.
- In Sun Solaris clients t is set adaptively in the range 30 to 60 seconds depending on file update frequency. The range is 30 to 60 seconds there is a lower risk of concurrent update.



- The validity check is made on each access to cache block. If the first half of the check is true, then the second half of the check need not be made. The first half of the check does not require network I/O.
- A separate first half of the check for the client retrieves $T_{m, server}$ valid, and the client
- If they do not match, then the cache block is invalid, and the client must request a new copy from the server.
- Traffic can be reduced by applying new value blocks and by piggy-backing attribute values
- Write-back is used for writes, where modified files are flushed when a file is closed or when a sync operation takes place in the VFS. Special purpose daemons are used to do this asynchronously.



- **Access transparency** : Yes. Applications programs are usually not aware that files are remote and no changes are need to applications in order to access remote files.
- **Location transparency** : Not enforced. NFS does not enforce a global namespace since client file systems may mount shared file systems at different points. Thus an application that works on one client may not work on another.
- **Mobility** : Not supported for updates. A file must be updated more than once.
- **Scalability** : Not supported for updates. A file must be updated more than once.
- **File replication** : Not supported for updates. A file must be updated more than once.
- **Hardware and operating system heterogeneity** : Almost every known operating system and hardware platform can run NFS.
- **Fault tolerance** : Acceptable. NFS is stateless and does not require a server to handle failures.
- **Consistency** : Tunable. NFS is not recommended for close synchronization between processes.
- **Security** : Kerberos is integrated with NFS. Secure RPC is also an option being developed.
- **Efficiency** : Acceptable. Many options exist for tuning NFS.