

Lecture 25: April 27

*Lecturer: Prashant Shenoy**TA: Vimal Mathew & Tim Wood*

25.1 Distributed File Systems

A distributed file system is one of the most common types of distributed system. It is very common for a large number of computers to want to be able to have shared access to some set of data or program files. A distributed file system provides this by giving the abstraction that a shared disk appears to be attached to each node in the system, even though in reality it is only physically attached to one or more server nodes. This is commonly used in educational computer labs where users want access to their files regardless of which specific computer they log in to.

25.1.1 Naming and Transparency

In order to be useful, users must be able to address (or **name**) the shared files they would like to access. Different distributed file systems support naming in different ways. Some systems provide **location transparency**, which states that the knowing the name of a file does not reveal its physical storage location. This can be desirable for security reasons if administrators do not want users to necessarily know exactly where different files are stored. This can also produce a simpler system because it means that users do not *need* to know the exact storage location, either. **Location independence** says that even if a file is moved between physical storage locations, its name will still remain valid. This can be beneficial since users do not need to repeatedly look up where a file is located even if it is moved between disk drives. In practice, many distributed file systems support location transparency, but few provide location independence.

25.1.2 Naming Strategies

An **absolute name** is one which provides a complete address to a file including both the server and path names. This has the advantage that it is trivial to find a file once the name is given since it contains complete information. This means that no additional state must be kept since each name is self contained, which can lead to greater scalability. However, users must know the complete name of the server and path to files they want to access and because of the naming convention there is a clear difference between local and remote files. This limits transparency since users must deal with local and remote files differently, and makes it harder to changes such as when a file moves between hosts. This can also make the system less resilient to failure since there is no abstraction layer which could re-map addresses (names) if one server failed but another contained the same data. This technique is used in operating systems such as Windows and Apple.

An alternative approach is to use **mount points**. In this case, the client machine creates a set of “local names” which are used to refer to remote locations. These names are called mount points because the remote files or directories are connected to and attached to the local file system. The operating system must maintain a table (such as `/etc/fstab` in Linux) to maintain the mapping of what server and path are mapped to each mount point. When the system boots up, it scans through the table and connects to each remote server. It then translates any file accesses to the mount point into network calls which are transferred to the

remote file system. This provides location transparency described previously, because once the mapping is made, the local client does not need to know or care that the files it are using are actually located across the network. This allows the remote server to be changed without affecting the local file name mappings, although the system may need to be restarted to acknowledge the change. The main disadvantage of using mount points is that it can lead to confusion since two different local names may actually map to the same file on a remote system. This approach is commonly used in Linux and Sun operating systems.

A third option for naming is to use a **global name space**. In this case, all nodes within the system have an identical name space—the path and name of a file on one machine will be the same on every other machine, regardless of where the file is actually stored. This is typically implemented using a set of dedicated file servers that store all files for the system. When a client boots up, it contacts one of the file servers and receives the layout of the distributed file system. When a user accesses a file, the server sends a copy of it to the client machine where it is cached. As the client updates the files, the changes must be written back to the central file servers. The advantage of this approach is that naming is consistent across all clients. Also, the storage servers are able to seamlessly move files around because clients always contact the server to learn where files are located within the global name space. However, the fact that files are cached by clients can lead to challenges in keeping file content consistent across all nodes. Enforcing a global name space across all nodes also limits flexibility, and can lead to performance problems, particularly when the scale of the system grows.

25.1.3 Remote File Access and Caching

When a client wants to modify a file contained in a distributed file system, it can perform all operations remotely, or use local caching. In the first case, the client sends the server an operation such as a write call, and the server performs the write and returns a result (often this is done using RPC). Alternatively, when a write is to be performed, the server can send the file (or part of the file) back to the client so that it can make the changes locally; this results in the client building up a local cache of files which it is working with.

If caching is used, then the distributed system faces additional challenges since it must figure out how to propagate changes back to the server after they have been made. This is especially complicated when multiple clients may be accessing the same file—the system must preserve consistency so that one clients actions do not conflict with those made by another.

A second challenge is deciding whether clients should keep its cache in memory or on its local disk. Using the local disk is significantly faster than making accesses over the network, but it is still much slower than keeping the file in memory. However, the disk is a persistent store, so it can provide greater reliability in the case that the node fails before writing its changes back to the file server. Of course, writing to disk also requires that the client *has* a disk, which is not always the case in “thin client” systems. Alternatively, the files can be cached in memory, but while this provides greater speed, it is less reliable and results in a more limited cache size.

25.1.4 Cache Update Policy

A cache update policy defines when writes made to a cache should be propagated back to the original disk. Using a **write through** policy provides high reliability since writes are immediately written to the server. However, the user will see reduced performance because all writes not only need to be made to the local disk, but transmitted over the network and acknowledged by the server. As a result, the cache provides no benefit for write requests (identical to doing all writes remotely), but it can still improve performance for reads.

In a **write back** policy, writes to the remote server are delayed until some event occurs such as the file being closed, the block being removed from the cache, or some set delay. This leads to better performance because writes can be queued up and done asynchronously as the local program continues executing. This can also reduce network traffic by exploiting the fact that a single block may be written and overwritten multiple times in a short window—with a write back policy, these changes can be aggregated into a single write which is sent to the server. Unfortunately, write back cannot provide the same reliability guarantees as write through since the node may crash or lose power before the write is sent back to the file server.

Cache Consistency is another important issue that determines how caches are maintained when multiple clients may be accessing the same file. In a **client-initiated consistency** system, the client is responsible for checking with the server to verify that each file in its cache is consistent (e.g. that no one else has modified the file since it was cached). Depending on the level of consistency required, the client could verify that its cache is consistent on every access, at a given interval, or only when the file is first opened. This is a relatively simple protocol to implement, but it requires the servers to trust that clients will indeed verify that their caches are correct—a single corrupt or malicious client could disrupt the complete system.

In **server-initiated consistency**, the server acts as a central authority over which clients have up to date or invalid caches. In this case, the server must track information about all clients that have cached a file so it knows which parts of which files are currently cached, as well as whether a given client is reading or writing to a file. Using this information, the server is able to detect when reading and writing clients might conflict with each other, and will send messages to clients to force them to invalidate their cache entries and request them again.

25.1.5 Server State and Replication

When designing server software for a distributed system of any sort, the programmer must decide whether the server should be stateful or stateless. A **stateful** server is one which maintains information about the clients that are connected with it. This provides better performance because clients do not need to repeatedly inform the server of their identity of their past requests, reducing overhead. However, a stateful server can be less resilient to failures—if the server crashes it loses all of its state leading to errors or forcing clients to somehow send information to rebuild the state. A **stateless** server is one that does not maintain any state at all about client connections. In this case, every request made by a client must provide sufficient information for the server to authenticate the client and return a proper result. Stateless servers can be more fault tolerant because if a server crashes it can simply be rebooted—since it does not have any state to lose it will have no impact on the system other than a delay while the machine or process is restarted.

The distributed systems we have described can involve a large number of clients accessing a server. As a result, the server can easily become a bottleneck, reducing performance. **Replication** involves running multiple servers that each contain the data or services provided by the distributed system. This can improve both performance and fault tolerance since client requests can be distributed across multiple servers, and if one server fails the others can continue working. However, replicating servers adds difficulties similar to the cache consistency ones listed previously, since now the distributed system must ensure that all servers maintain consistent data.

See the slides for a description of how Sun's Network File System deals with many of these challenges.