



- (51) **International Patent Classification:**  
**H04L 12/24** (2006.01) **H04L 29/06** (2006.01)  
**G06F 15/16** (2006.01)
- (21) **International Application Number:**  
PCT/US2011/063618
- (22) **International Filing Date:**  
6 December 2011 (06.12.2011)
- (25) **Filing Language:** English
- (26) **Publication Language:** English
- (30) **Priority Data:**  
12/964,749 10 December 2010 (10.12.2010) US
- (71) **Applicant (for all designated States except US):** **MICROSOFT CORPORATION** [US/US]; One Microsoft Way, Redmond, Washington 98052-6399 (US).
- (72) **Inventors:** **SWAN, Paul R.**; c/o Microsoft Corporation, LCA - International Patents, One Microsoft Way, Redmond, Washington 98052-6399 (US). **GEORGE, Mathew**; c/o Microsoft Corporation, LCA - International Patents, One Microsoft Way, Redmond, Washington 98052-6399 (US). **KRUSE, David M.**; c/o Microsoft Corporation, LCA - International Patents, One Microsoft Way, Redmond, Washington 98052-6399 (US). **BATTEPATI, Roopesh C.**; c/o Microsoft Corporation, LCA - International Patents, One Microsoft Way, Redmond, Washington 98052-6399 (US). **JOHNSON, Michael C.**; c/o Microsoft

Corporation, LCA - International Patents, One Microsoft Way, Redmond, Washington 98052-6399 (US).

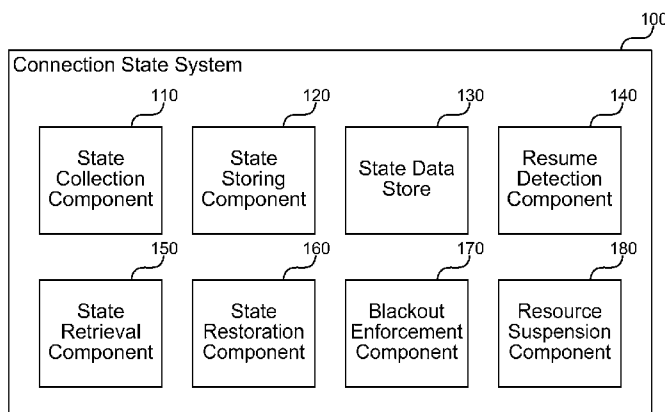
- (81) **Designated States** (unless otherwise indicated, for every kind of national protection available): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.
- (84) **Designated States** (unless otherwise indicated, for every kind of regional protection available): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

**Declarations under Rule 4.17:**

- as to applicant's entitlement to apply for and be granted a patent (Rule 4.17(ii))

[Continued on next page]

- (54) **Title:** PROVIDING TRANSPARENT FAILOVER IN A FILE SYSTEM



**FIG. 1**

- (57) **Abstract:** A connection state system is described herein that allows a client to resume a connection with a server or a different replacement server by remotely storing client state information in association with a resume key. The system provides a resume key filter operating at the server that facilitates the storing of volatile server state information. The state information can include information such as oplocks, leases granted to a client, and in-flight operations on a file handle. The resume key filter driver sits above the file system, which allows multiple file access protocols to use the filter. Upon a failover event, such as a server going down or losing connectivity to a client, the system can bring up another server or the same server and reestablish state for file handles held by various clients using the resume key filter.





- 
- *as to the applicant's entitlement to claim the priority of the earlier application (Rule 4.17(iii))*
- Published:**
- *without international search report and to be republished upon receipt of that report (Rule 48.2(g))*

## PROVIDING TRANSPARENT FAILOVER IN A FILE SYSTEM

### BACKGROUND

[0001] A variety of techniques exists for sharing files, printers, and other resources between two computers on a network. For example, two application-layer network protocols for sharing resources are Server Message Block (SMB) and Network File System (NFS). SMB is used by MICROSOFT™ WINDOWS™ and other operating systems to allow two computers or other resources to communicate, request access to resources, specify intended access of resources (e.g., reading, writing, etc.), lock resources, and so on. MICROSOFT™ WINDOWS™ Vista introduced SMB 2.0, which simplified the command set of SMB 1.0 and added many other enhancements. MICROSOFT™ WINDOWS™ 7 and Server 2008 R2 introduced SMB 2.1, which added opportunistic locking (oplocks) and other enhancements.

[0002] Most protocols for remote sharing of resources assume a one-to-one relationship between connections and sessions. A session represents the lifetime of any single request to access a resource and the subsequent access of that resource until the connection is terminated. A session may also be associated with a particular security principal and validated security credentials that determine the actions that are authorized during the session. A connection can include a Transmission Control Protocol (TCP), User Datagram Protocol (UDP), or other type of connection over which higher-level protocols like SMB and NFS can communicate to carry out commands. An SMB or NFS session typically involves opening a TCP or UDP connection between a source of a request and a target of the request, sending one or more SMB or NFS commands to access the target resource, and then closing the session. Sometimes connections are lost during a session (e.g., due to a network failure), tearing down any client and server state established during the connection. To reestablish a connection the client and server typically have to repeat all of the steps used to initially establish the connection over again.

[0003] The SMB2 protocol provides a resume key that allows clients to quickly reestablish a file handle to a server if a client is disconnected from the server, enabling clients to reduce network round trips to the server and reduce the load on the server when a client reconnects. However, today the resume key does not provide restoration of state in the event of server failover in which the SMB2 server loses volatile state during a server reboot or failover of a cluster. State information associated with existing opens is

lost and must be reestablished. In addition, the resume key is an application-level concept that can only be created and used within the boundary of an application but not shared.

#### SUMMARY

[0004] A connection state system is described herein that allows a client to resume a connection with a server or a different replacement server by remotely storing client state information in association with a resume key. The system provides a resume key filter operating at the server that facilitates the storing of volatile server state information. The state information can include information such as oplocks, leases granted to a client, and in-flight operations on a file handle. The resume key filter driver sits above the file system, which allows multiple file access protocols to use the filter, as well as permitting the filter to provide this functionality across multiple file systems. The system provides state information to the protocol, independent of the actual protocol. Upon a failover event, such as a server going down or losing connectivity to a client, the system can bring up another server or the same server and reestablish state for file handles held by various clients using the resume key filter. The filter enforces a blackout window on active files after failover that guarantees that the active file state can be consistently restored and that other clients do not step in to access the file in the interim. In the resume phase, the resume key is used to map existing pre-failover file handles to post-failover preserved file state stored by the resume key filter. Thus, the connection state system allows the same or another server to resume the state of a previous session with a client after a failover event with as little disruption as possible to clients.

[0005] This Summary is provided to introduce a selection of concepts in a simplified form that are further described below in the Detailed Description. This Summary is not intended to identify key features or essential features of the claimed subject matter, nor is it intended to be used to limit the scope of the claimed subject matter.

#### BRIEF DESCRIPTION OF THE DRAWINGS

[0006] Figure 1 is a block diagram that illustrates components of the connection state system, in one embodiment.

[0007] Figure 2 is a flow diagram that illustrates processing of the connection state system to capture file system state information, in one embodiment.

[0008] Figure 3 is a flow diagram that illustrates processing of the connection state system to resume a connection after failover, in one embodiment.

[0009] Figure 4 is a block diagram that illustrates the operating environment of the connection state system, in one embodiment.

## DETAILED DESCRIPTION

**[0010]** A connection state system is described herein that allows a client to resume a connection with a server or a different replacement server by remotely storing client state information in association with a resume key. The system provides a resume key filter  
5 operating at the server that facilitates the storing of volatile server state information. The state information can include information such as oplocks, leases granted to a client, and in-flight operations on a file handle. The resume key filter driver sits above the file system, which allows multiple file access protocols to use the filter, as well as permitting the filter to provide this functionality across multiple file systems. The system provides  
10 state information to the protocol, independent of the actual protocol. Upon a failover event, such as a server going down or losing connectivity to a client, the system can bring up another server or the same server (e.g., via different connection, such as a redundant Ethernet connection) and reestablish state for file handles held by various clients using the resume key filter.

**[0011]** The system provides a resume key filter that can be used for transparent failover after a server loses its connection to a client. The resume key filter sits atop the file system and is therefore independent of protocol used to access the file system. The resume key filter records active file state and then restores the active file state after a failover. The resume key filter can capture a variety of state information. For example,  
20 the filter records the active file system state comprising open handles (statically referenced by a resume key), uncommitted file state (such as delete on close, delete pending, and lock state), and certain in-flight/interrupted file operations. The filter restores the active file system state after failover such that the open handles are resumed to match those prior to failover and in-flight operations can be consistently replayed. The filter provides a means  
25 for multiple Remote File Systems (RFS) to store and retrieve private opaque data that is associated with an open file handle referenced through a resume key. The filter enforces a blackout window on active files after failover that guarantees that the active file state can be consistently restored and that other clients do not step in to access the file in the interim. The filter also allows a currently active file to be “suspended” and then resumed  
30 without a failover in order to support SMB in the cluster scenario where nodes failover.

**[0012]** A remote file system (RFS) supplies a resume key with every file create operation as an extra parameter during create. The key is unique to the RFS. The resume key filter uses a resume key and an RFS identification key together as a globally unique identifier (GUID) for a file handle. In the resume phase, the resume key is used to map

existing pre-failover file handles to post-failover preserved file state stored by the resume key filter. Thus, the connection state system allows the same or another server to resume the state of a previous session with a client after a failover event with as little disruption as possible to clients.

5 [0013] Figure 1 is a block diagram that illustrates components of the connection state system, in one embodiment. The system 100 includes a state collection component 110, a state storing component 120, a state data store 130, a resume detection component 140, a state retrieval component 150, a state restoration component 160, a blackout enforcement component 170, and a resource suspension component 180. Each of these components is  
10 described in further detail herein.

[0014] The state collection component 110 creates a state record for each file handle and collects state information as a client request operations using the file handle. The component 110 may operate at a server and store state information externally from the server so that the state information can be accessed if the server is unavailable. For  
15 example, the component 110 may store the state information the state data store 130 described further herein. The state collection component 110 may receive a resume key from the client when the client connects to the server, and the component 110 associates collected state information with the resume key in the state data store 130. If a client is reconnecting after a failover event, the client will provide the same resume key used to  
20 open the initial connection and the current server can find the state information stored by the previous server and recreate the server state from the state information.

[0015] The state storing component 120 stores collected state information in association with a resume key provided by the client. The component 120 stores the state information in the state data store 130 and keeps a record of operations related to the resume key that  
25 would be restored in the event of a failover event. The state information may include open file handles, oplocks granted, leases and lease information, in-progress file operations, byte range locks, and any other information that another server would use to carry out the client's requests without the client reestablishing all of the previous state.

[0016] The state data store 130 persistently stores file system state information that a  
30 resuming server uses to recreate state information stored by a failing server. In some cases, the resuming server and the failing server may be the same server using a different connection to the client or coming back up after a brief outage. In other cases, the resuming server and failing server are different servers, and the state data store 130 is provided in a location accessible to both servers for sharing the state information. The

state data store 130 may include one or more files, file systems, hard drives, databases, storage area networks (SANs), cloud-based storage services, or other storage facility for persistently storing data and accessible to both the failing and resuming server for exchanging information. As the failing server is performing operations, it is storing state information about the operations' progress in the state data store 130. Upon a failure, the failing server will be interrupted, and a resuming server accesses the state information to resume the state and continue carrying out any operations that did not complete.

[0017] The resume detection component 140 detects a condition that makes a failing server unavailable and informs a resuming server to act in the failing server's place. The detection may be client driven, such that the system does not perform any resuming steps until the client reconnects to the system and provides a previously used resume key. The system identifies the key and any state information stored in association with the key and restores that state information as part of setting up the connection. The resuming server may be the same or a different server from the failing server, and the resume detection component 140 ensures that the resuming server becomes active to handle the client's requests. In other embodiments, the detection may be server driven and the system may proactively bring up a resuming server upon detecting that a failing server has gone down. The system may also prepopulate the resuming server with stored state information even before a client requests a connection to the server.

[0018] The state retrieval component 150 retrieves stored state information from a location accessible to the resuming server, wherein the state information allows the resuming server to resume any previously requested file system operations that were interrupted by the detected failure condition. The state retrieval component 150 retrieves state information from the state data store 130 and invokes the state restoration component 160 to load the information into the resuming server so that the resuming server can continue the operations requested by the client.

[0019] The state restoration component 160 loads the retrieved state information into the resuming server so that the resuming server can continue operations previously requested by the client. The restoration may also include refreshing any oplocks and/or leases held by the client to ensure that other clients abide by previously requested access levels and/or exclusivity granted to the client. The state restoration component 160 allows a new server or node to take the place of a failing server or node without placing a heavy burden upon the client to restore state information by repeating past operations. Clients using protocols like SMB 2.0 already know how to use a resume key to restore a connection to the same

server, and the connection state system allows a substitute server to take the place of a failing server transparently to the client. Resume keys can also be used with NFS. In the case of NFS, the concept of a resume key is completely opaque to the client. The client does not explicitly refer to or participate in resume key generation, management, and association. Rather, the resume key is a server side concept.

**[0020]** The blackout enforcement component 170 enforces a blackout period on access to one or more files or other resources that prevents a second client from interfering with resources in a way that would conflict with a first client resuming a connection to the resuming server. The component 170 may automatically select a period deemed to be long enough to avoid most conflicting operations (e.g., 15 or 30 seconds), but not so long as to prevent other clients from accessing resources if the first client does not resume the connection. The period allows the first client time to resume the connection if the first client chooses. In some embodiments, the system allows an administrator or other user to configure the duration of the blackout period to tune the system for application-specific purposes. The system may also allow individual clients to request a blackout period as a parameter to a create/open request or other application programming interface (API). In response to attempts to access a blacked out resource, the component 170 may provide an indication to try again after a particular period or simply fail the request. After the blackout period if no client has resumed the connection, then the blackout is over and requests to access the resource will succeed as normal.

**[0021]** The resource suspension component 180 allows a currently active resource to be suspended and resumed without a failover event to allow a cluster to failover to another node in planned manner. One example is load balancing. Suspending allows scenarios where a subset of the state is being transitioned to a new node. For example, if one node in the cluster is overloaded, an administrator may want to migrate half the node's clients to a new node. Suspending allows capturing the state of the opens that are being migrated and allows the client to connect to the new node as a continuation of the same open (e.g., without reestablishing server state). As another example, SMB supports clustering scenarios in which generic nodes are brought into a cluster and can be used interchangeably to service client requests. Sometimes there is a reason to bring down a particular node, such as for maintenance, and it is desirable to cleanly suspend the current node, activate the new node, deactivate the old node, and then perform any maintenance operations on the deactivated node. This can have an undesirable impact on clients, but



using the techniques described herein, the system 100 can suspend the node in an organized manner, and allow clients to resume operations with the new node efficiently.

[0022] The computing device on which the connection state system is implemented may include a central processing unit, memory, input devices (e.g., keyboard and pointing devices), output devices (e.g., display devices), and storage devices (e.g., disk drives or other non-volatile storage media). The memory and storage devices are computer-readable storage media that may be encoded with computer-executable instructions (e.g., software) that implement or enable the system. In addition, the data structures and message structures may be stored or transmitted via a data transmission medium, such as a signal on a communication link. Various communication links may be used, such as the Internet, a local area network, a wide area network, a point-to-point dial-up connection, a cell phone network, and so on.

[0023] Embodiments of the system may be implemented in various operating environments that include personal computers, server computers, handheld or laptop devices, multiprocessor systems, microprocessor-based systems, programmable consumer electronics, digital cameras, network PCs, minicomputers, mainframe computers, distributed computing environments that include any of the above systems or devices, set top boxes, systems on a chip (SOCs), and so on. The computer systems may be cell phones, personal digital assistants, smart phones, personal computers, programmable consumer electronics, digital cameras, and so on.

[0024] The system may be described in the general context of computer-executable instructions, such as program modules, executed by one or more computers or other devices. Generally, program modules include routines, programs, objects, components, data structures, and so on that perform particular tasks or implement particular abstract data types. Typically, the functionality of the program modules may be combined or distributed as desired in various embodiments.

[0025] Figure 2 is a flow diagram that illustrates processing of the connection state system to capture file system state information, in one embodiment. Beginning in block 210, the system receives from a client a request to access a remote resource stored on a server. The access request may include one or more parameters, including a resume key used to identify the session across multiple potential connections if a connection fails. The resource access request may be the first in a series of access requests sent from the client, and if the client is ever disconnected from the server the client may provide the same resume key in a subsequent open request to the same or a new server to resume the

connection. The resume key helps the server respond to the client faster by correlating state information maintained by the server (or across servers) between what would otherwise appear to be independent client connections.

[0026] Continuing in block 220, the system determines an identifier that identifies a client session related to the request. The identifier in some cases is a resume key that the client provides for durable handles that allow resuming sessions that get disconnected for various reasons. The access request may include one or more parameters at well-defined locations in the protocol so that the system can extract the key by reading the appropriate location in the request. Alternatively or additionally, the server may include an automated process for determining the identifier that does not involve information explicitly provided by the client. For example, the server may identify the client by Internet Protocol (IP) address or other inferred data that indicates to the server that the client connection is correlated with a previous session.

[0027] Continuing in block 230, the system creates a resume record searchable by the extracted identifier that associates state information created by operations requested by the client with the extracted identifier. The resume record can be stored at a location external to the server handling the present access request so that if the server fails another server will be able to read the record to resume the operations and act in the original server's place. The resume record may include a file, database record, or other form of storage. The record may contain a list of open file handles, oplocks obtained by the client, leases, or other file system state information.

[0028] Continuing in block 240, the system receives a file operation from the client that requests access to a file accessible through the server. The file operation may be a request to open a file, close a file, read a file, write a file, print to a shared printer, or other file system operations. The received operation involves a certain amount of state information being created on the server. For example, if the client opens a handle to the file, then the server tracks that handle to manage other client requests related to the file and to manage lifetime and/or cleanup processing for the handle.

[0029] Continuing in block 250, the system stores resume state information in the created resume record that provides information to resume the received file operation if the client loses its connection with the server. If the client connection fails, the client will attempt to resume the connection by again opening a remote resource and specifying the same resume key or other session identifier. This will allow the server or another server to access the stored resume record and reestablish the previous state information.

[0030] Continuing in block 260, the system performs the requested file operation. The operation may open a file, read the contents of the file, write data to the file, change access rights to the file, or any other file system operation. The outcome of the operation may change the state stored by the server. For example, if the client attempts to close a handle and the server successfully closes the handle, then the server state will be updated to remove the handle from a list of handles tracked by the server.

[0031] Continuing in block 270, the system updates the stored resume state information in the created resume record based on an outcome of the performed file operation. The system cannot know in advance when a failure will occur that causes failover, so the system keeps an up to date view of the server state in the resume record that allows a server to reestablish the state as close to the previous server's state as possible.

Operations that were not completed may be replayed to complete the operations while operations that did complete will not need to be repeated (but the server may resend the result to the client). Thus, the system updates the state as needed during and after various file system operations that change server state information.

[0032] Continuing in block 280, the system sends a response to the client that indicates the outcome of the requested file operation. If the client and server are still connected, then operations continue as they are requested by the client and the server continues to track updated state information. If at any time the connection is lost, another server can be brought up or the existing server repaired and the state information can be loaded from the state store to reestablish the prior server state. Upon receiving a new request from the client to resume the session, the client need not be aware that failover has occurred and that the client is potentially interacting with a different server than the original one. After block 280, these steps conclude.

[0033] Figure 3 is a flow diagram that illustrates processing of the connection state system to resume a connection after failover, in one embodiment. Beginning in block 310, the system receives from a client a request to open a remote resource stored on a server. The access request may include one or more parameters, including a resume key used to identify the session across multiple potential connections if a connection fails. Unlike the resource access request discussed with reference to Figure 2, this request is a request to reconnect to a previously connected session. The client provides the same resume key as originally provided, so that the server can correlate the current session request with the previous session.

[0034] Continuing in block 320, the system determines a session identifier that identifies a client session related to the request. The identifier in some cases is an SMB 2 resume key that the client provides for durable handles that allow resuming sessions that get disconnected for various reasons. The access request may include one or more parameters at well-defined locations in the protocol so that the system can extract the key by reading the appropriate location in the request. In other cases, the server may determine the identifier automatically based on information about the client.

[0035] Continuing in block 330, the system looks up the received session identifier in a state store to identify a resume record associated with the session identifier. Any previous server interacting with the client using a resumable session stores state information on an ongoing basis throughout interaction with the client. When the client attempts to reestablish the connection, the state information is available to a failover server standing in for the original server. The state information may be stored externally to the original server so that the information is accessible after a failure of the original server.

[0036] Continuing in block 340, the system receives from the state store previous state information associated with the resume record. The state information identifies static state, such as open file handles, obtained leases, obtained oplocks, and so forth, as well as dynamic state, such as in-flight operations that may not have completed. The stored state information allows the failover server to take the place of the original server without specific processing by the client. The client understands resumable handles and performs steps to make a connection resumable, but may not be aware of which server ends up handling the connection at any particular time. The client may access the server via a domain name or network file share that can resolve to an address of any one of several servers, including the failover server.

[0037] Continuing in block 350, the system restores the received previous state information by loading the information into the file system components that track file system state. After loading the state, the local state of the failover server is similar to how the state would look if all of the previous operations had occurred on the failover server. Thus, the failover server is as useful to the client for continuing the series of operations as the original server would have been had the connection not failed.

[0038] Continuing in block 360, the system responds to the client access request indicating that the server found the resume record and is ready to receive client operations related to the previous session. Based on the server's response, the client can determine whether the session is resumed or whether the client needs to take steps to repeat previous

operations. If the session was successfully resumed, then the client can continue knowing the previous operations completed or were replayed to complete after the server resumed. In some cases, the system may hand the client a new file handle that has the same state as the pre-failover file handle. After block 360, these steps conclude.

5 [0039] Figure 4 is a block diagram that illustrates the operating environment of the connection state system, in one embodiment. The environment includes one or more operating system services or applications that interact with file systems. For example, MICROSOFT™ WINDOWS™ includes a server service 420 known as SRV, and a network file system service 410 known as NFS. The network file system service 410 and  
10 server service 420 provides access to shared resources, such as files and printers, between computer systems. The server service 420 uses the SMB protocol common to WINDOWS™ networks, while the network file system service 410 provides access to Unix-based systems that more commonly use NFS. Regardless of the protocol, the resume key filter 430 captures file operations and stores state information for resuming the  
15 operations in a remote data store. The operations pass through the file system level 440 (e.g., NTFS or other file system), and affect one or more user data files 450. Meanwhile, the resume key filter 430 writes state information to a log file 460 or other data store, that another server can access to retrieve state information and resume a connection to a client. The system can operate independent of the particular protocol or file system involved, and  
20 various components can be updated to save their own particular state information in the state data store.

[0040] In some embodiments, the connection state system stores opaque blobs of data on behalf of file system components to allow the system to resume connections without component-specific knowledge. For example, the resume key filter described herein can  
25 ask the server service for any data the server service would need to recreate its present state. The filter can then store any received data as an opaque blob (i.e., the filter need not know what is in the blob or its semantic meaning) in the state store. Upon a failover condition, a resume key filter operating on the new server can access the stored state information, retrieve the stored blob, and provide the blob to the server service so that the  
30 server service can restore its own state. In this way, the system can be made to work with many types of protocols without specific knowledge of the internal operations of components that implement each protocol for a server.

[0041] In some embodiments, the connection state system blocks other clients from accessing files or other resources related to a resumable handle for some amount of time

(i.e., blackout period). If the original client reconnects during the blackout period, then the original client gets its connection back with all of the previous state, and can resume operations. If another client attempts to connect, the server may provide a message indicating to wait an amount of time and retry. Resume aware clients can use this  
5 information to delay retrying until after the blackout period, while older clients may simply fail the connection and manually retry at the user's request. If the original client does not return within the blackout period, the server cleans up the resume state information and allows new clients to access the resources as usual.

[0042] In some embodiments, the connection state system can use a variety of storage  
10 devices or strategies for speeding up resumes. For example, the system may use a fast, nonvolatile storage device (e.g., a solid state disk (SSD)) for storing resume state information so that resumes get faster access to data to avoid delaying operations already interrupted by a failure any further. As another example, the system may broadcast all changes made by each server to a group of servers, so that each server can maintain its  
15 own copy of the state information and can be the elected failover server in the event of a failure of the original server.

[0043] From the foregoing, it will be appreciated that specific embodiments of the connection state system have been described herein for purposes of illustration, but that various modifications may be made without deviating from the spirit and scope of the  
20 invention. Accordingly, the invention is not limited except as by the appended claims.

## CLAIMS

1. A computer-implemented method for capturing file system state information to facilitate resuming connections, the method comprising:
  - receiving from a client a request to access a remote resource stored on a server;
  - 5 determining an identifier that identifies a client session related to the request;
  - creating a resume record searchable by the extracted identifier that associates state information created by operations requested by the client with the extracted identifier;
  - receiving a file operation from the client that requests access to a file accessible through the server;
  - 10 storing resume state information in the created resume record that provides information to resume the received file operation if the client loses its connection with the server;
  - performing the requested file operation;
  - updating the stored resume state information in the created resume record based on
  - 15 an outcome of the performed file operation; and
  - sending a response to the client that indicates the outcome of the requested file operation,
  - wherein the preceding steps are performed by at least one processor.
2. The method of claim 1 wherein the access request includes one or more
- 20 parameters, including a resume key that identifies the client session across multiple potential connections if a connection fails, and wherein the resume key is at least part of the determined identifier.
3. The method of claim 1 further comprising, upon the client becoming disconnected from the server, receiving at a failover server a new access request that that the server can
- 25 correlate with the original access request to help the failover server respond to the client faster after a connection failure by correlating state information maintained by the server between multiple client connections.
4. The method of claim 1 wherein a Network File System (NFS) server determines the identifier automatically without receiving a resume key from the client.
- 30 5. The method of claim 1 wherein the extracted identifier is a Server Message Block (SMB) resume key that the client provides for durable handles that allow resuming sessions that get disconnected.
6. The method of claim 1 wherein creating the resume record comprises storing the resume record at a location external to the server handling the present access request so

that if the server fails another server will be able to read the record to resume any operations from the client and act in the original server's place.

7. The method of claim 1 where receiving the file operation comprises a request to perform an operation selected from the group consisting of opening a file, closing a file,  
5 reading a file, writing a file, obtaining a lease on a file, and obtaining a lock on a file.

8. The method of claim 1 further comprising, upon the client becoming disconnected from the server, loading at a failover server the stored resume record so that the client can connect to the failover server and continue any previous operations.

9. The method of claim 1 wherein performing the requested file operation modifies  
10 state stored by the server, and wherein updating the stored resume state information captures the modified state.

10. The method of claim 1 wherein updating the stored resume state information comprises keeping an up to date view of the server state in the resume record that allows another server to reestablish the state and handle client requests in place of the original  
15 server without requiring the client to reestablish at least some of the state information.

11. A computer system for providing transparent failover for clients in a file system, the system comprising:

a processor and memory configured to execute software instructions embodied within the following components;

20 a state collection component that creates a state record for each file handle and collects state information as a client requests operations using the file handle;

a state storing component that stores collected state information in association with a session identifier provided by the client;

25 a state data store that persistently stores file system state information that a resuming server uses to recreate state information stored by a failing server;

a resume detection component that detects a condition that makes a failing server unavailable and informs a resuming server to act in the failing server's place;

30 a state retrieval component that retrieves stored state information from a location accessible to the resuming server, wherein the state information allows the resuming server to resume any previously requested file system operations that were interrupted by the detected failure condition; and

a state restoration component that loads the retrieved state information into the resuming server so that the resuming server can continue operations previously requested by the client.



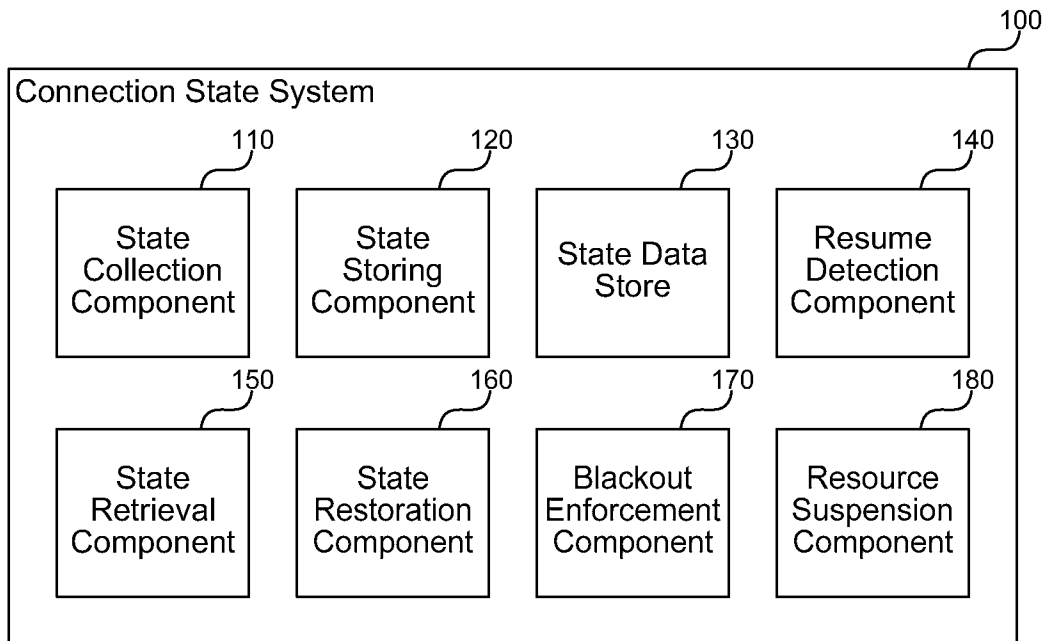
12. The system of claim 11 wherein the state collection component is further configured to operate at a server and store state information externally from the server so that the state information can be accessed if the server is unavailable.

13. The system of claim 11 wherein the state collection component is further  
5 configured to receive a resume key from the client when the client connects to the server, and associate collected state information with the resume key in the state data store.

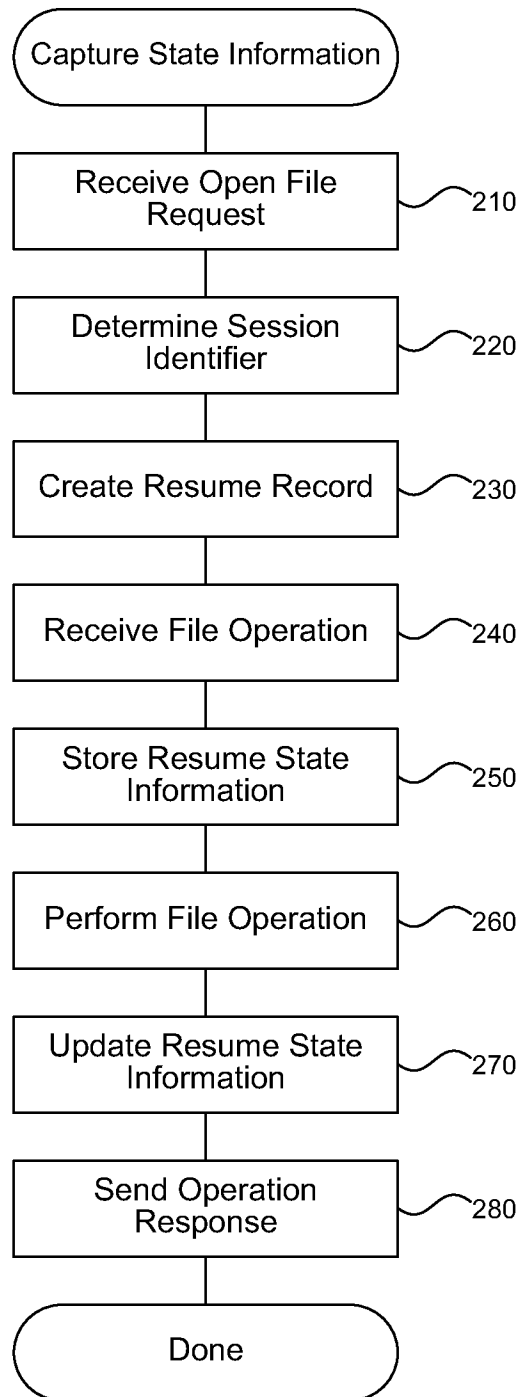
14. The system of claim 11 wherein the state data store stores and provides information for a resuming server that is the same server as the failing server using a different connection to the client.

10 15. The system of claim 11 wherein the state data store receives state information as the failing server is performing operations, and, upon a failure, provides access to the previously received state information to the resuming server to resume the state and continue carrying out any operations that did not complete.

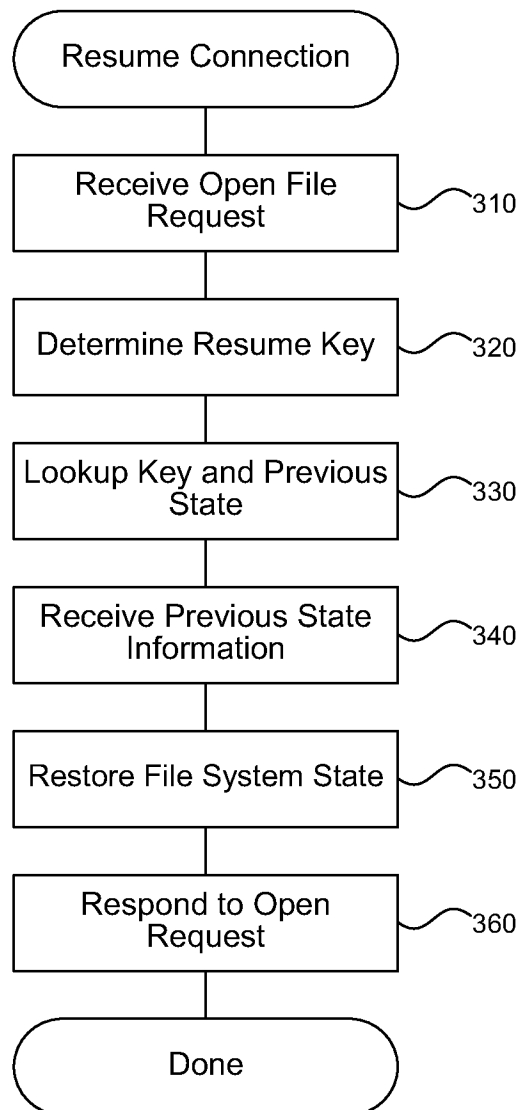
1/4

**FIG. 1**

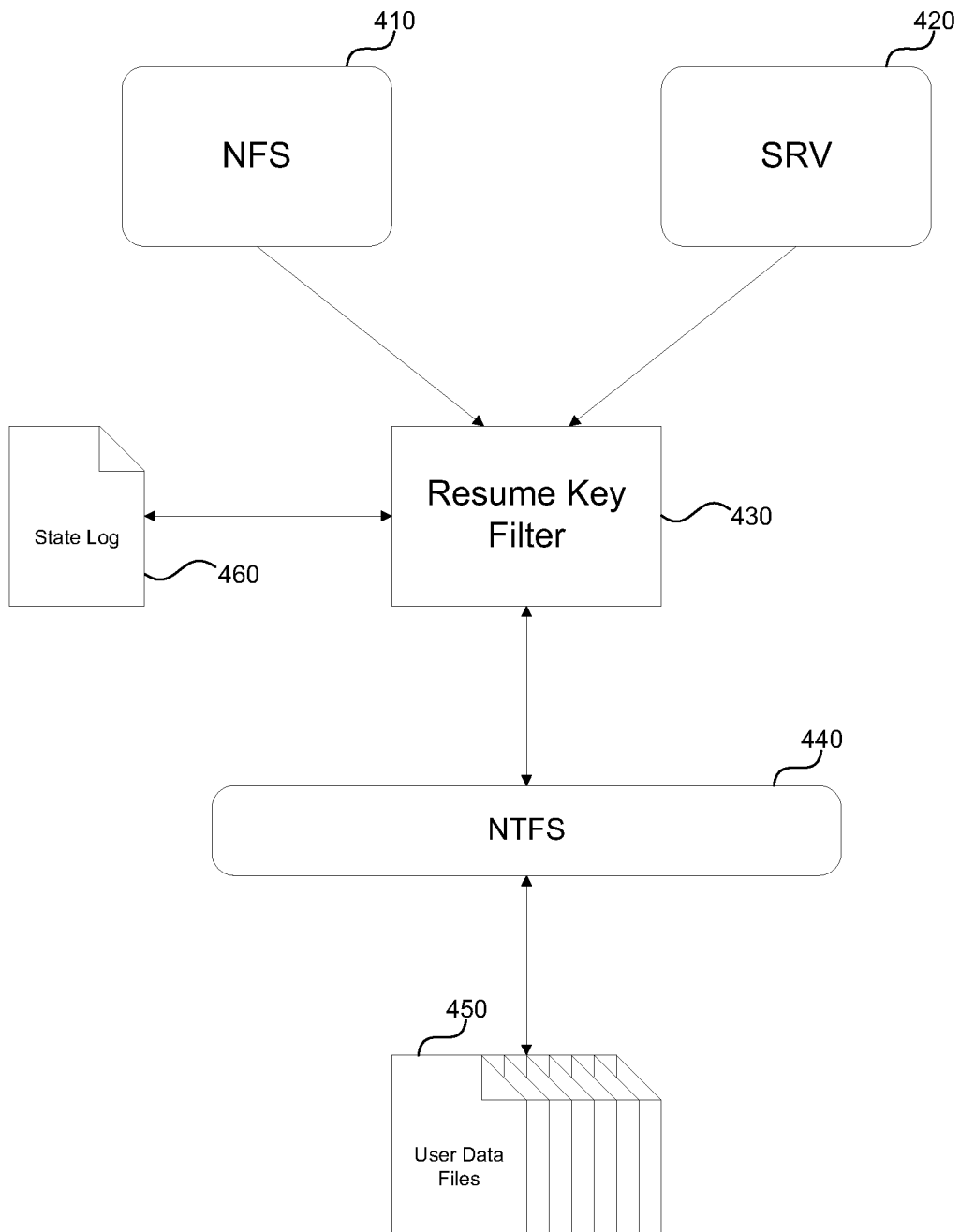
2/4

**FIG. 2**

3/4

**FIG. 3**

4/4

**FIG. 4**