

# CALLING THE WITNESS: SMB3 Failover with Samba/CTDB

**Günther Deschner**  
<gd@samba.org>  
Sr. Software Engineer  
Samba Team Member

**José A. Rivera**  
<jarrpa@samba.org>  
Software Engineer  
Samba Team Member

# About Samba, Red Hat, and Us

- **Currently 7 Samba Team members inside Red Hat**
- **Developers and users of Samba technology for authentication and storage solutions**
- **gd: 11 years Samba Team member  
8 years Red Hat (Samba Maintainer, Identity, Storage)**
- **jarrpa: 9 years working with Microsoft protocols  
3 years Red Hat (Samba Maintainer, Storage)**

# Agenda

- **Witness?**
- **Failover in SMB1/SMB2**
- **Failover in SMB1/SMB2 with CTDB**
- **Failover in SMB3**
- **The Witness Protocol**
- **Roadmap for Witness support in Samba**
- **Further reading & Q/A**

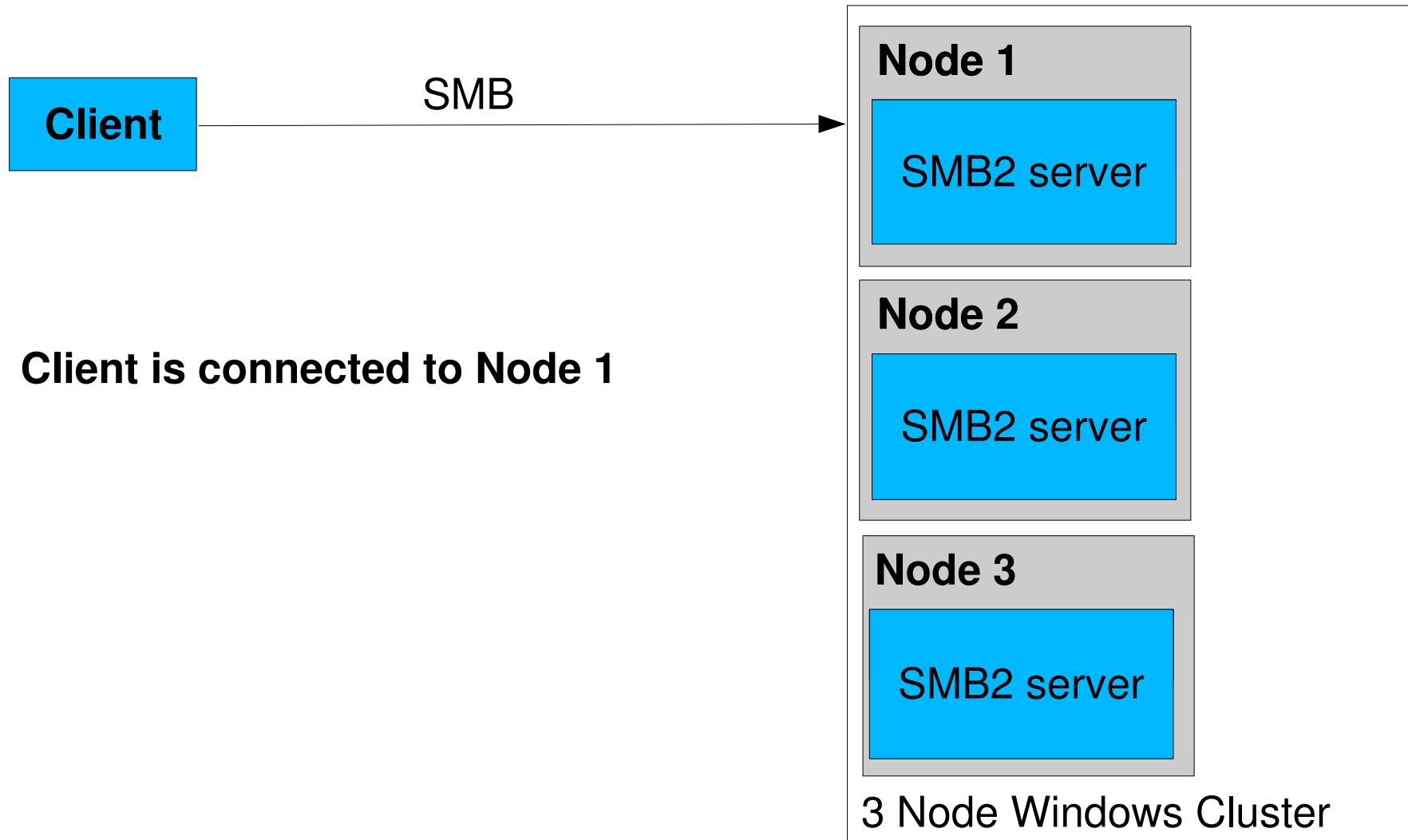
# Witness?

- **New DCE/RPC Service to “witness” availability of other services, in particular SMB3 connections**
- **Prompt and explicit notifications about failures in highly available systems**
- **Allows Continuous Availability of SMB shares in clustered environments**
- **Controlled way of dealing with reconnects instead of detecting failures due to timeouts**
- **Available with SMB3**

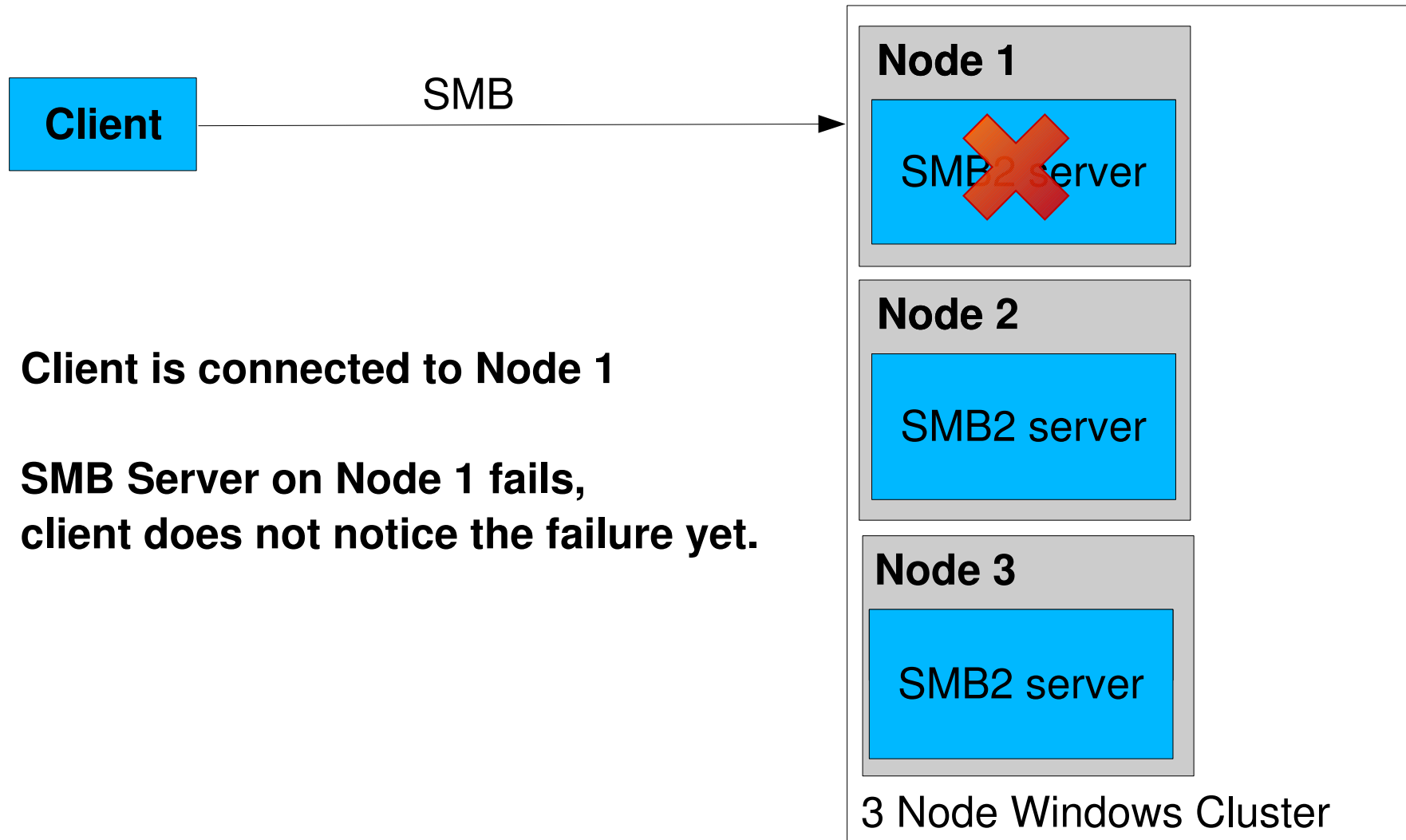
# Failover in SMB1/SMB2

- **Uncontrolled, clients detect unavailability by running into timeouts or by using keep alive mechanisms**
- **Clients reconnect after TCP/IP connection timeout**
- **Slow, unreliable, unpredictable**
- **Not all applications deal with stale connections good enough**

# Failover in SMB1/SMB2



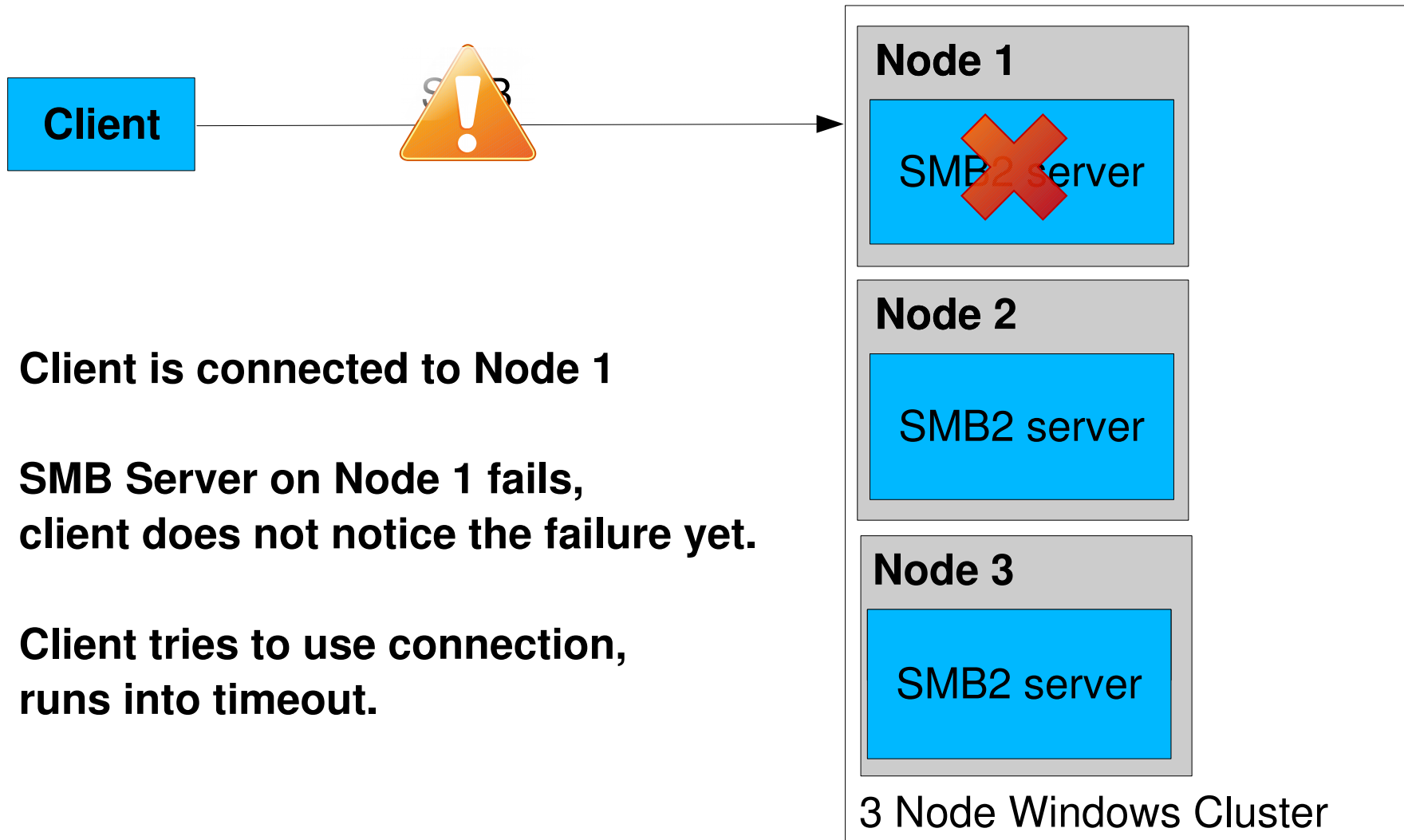
# Failover in SMB1/SMB2



**Client is connected to Node 1**

**SMB Server on Node 1 fails,  
client does not notice the failure yet.**

# Failover in SMB1/SMB2



**Client is connected to Node 1**

**SMB Server on Node 1 fails,  
client does not notice the failure yet.**

**Client tries to use connection,  
runs into timeout.**



# Failover in SMB1/SMB2

Client

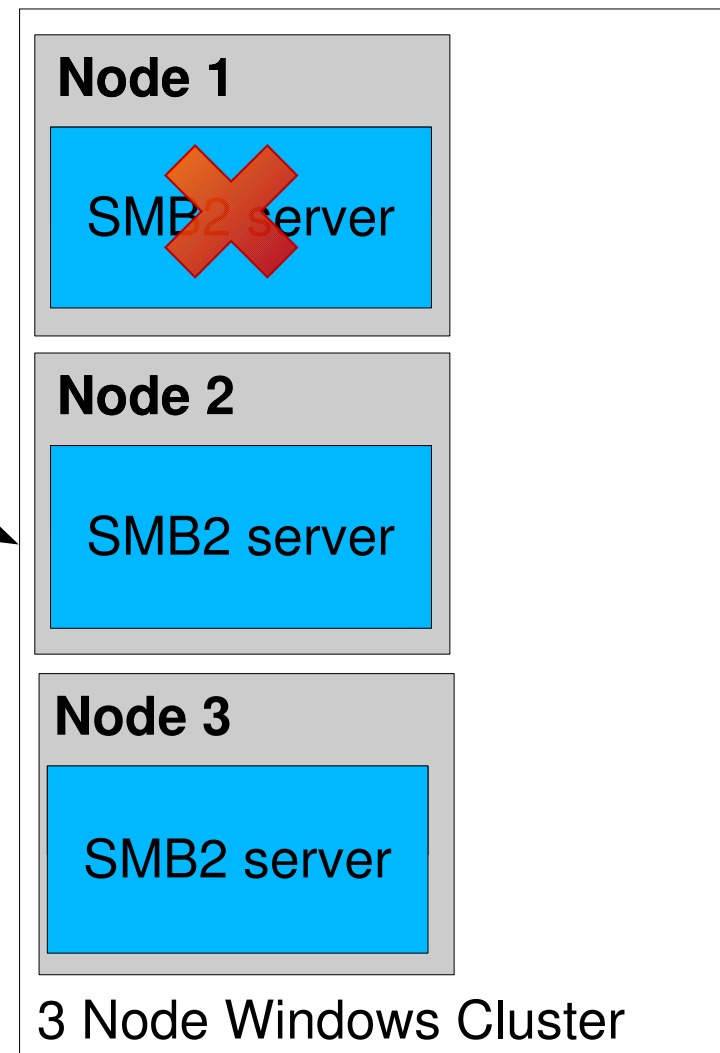
SMB

Client is connected to Node 1

SMB Server on Node 1 fails,  
client does not notice the failure yet.

Client tries to use connection,  
runs into timeout.

Finally Client reconnects to Node 2



# Failover in SMB1/SMB2 with CTDB

- Since 2007, a Samba cluster with CTDB is usually aware of failures before the client is
- In case of failure CTDB can proactively route the clients to another node
- With CTDB the cluster coordinates the failover, not the client

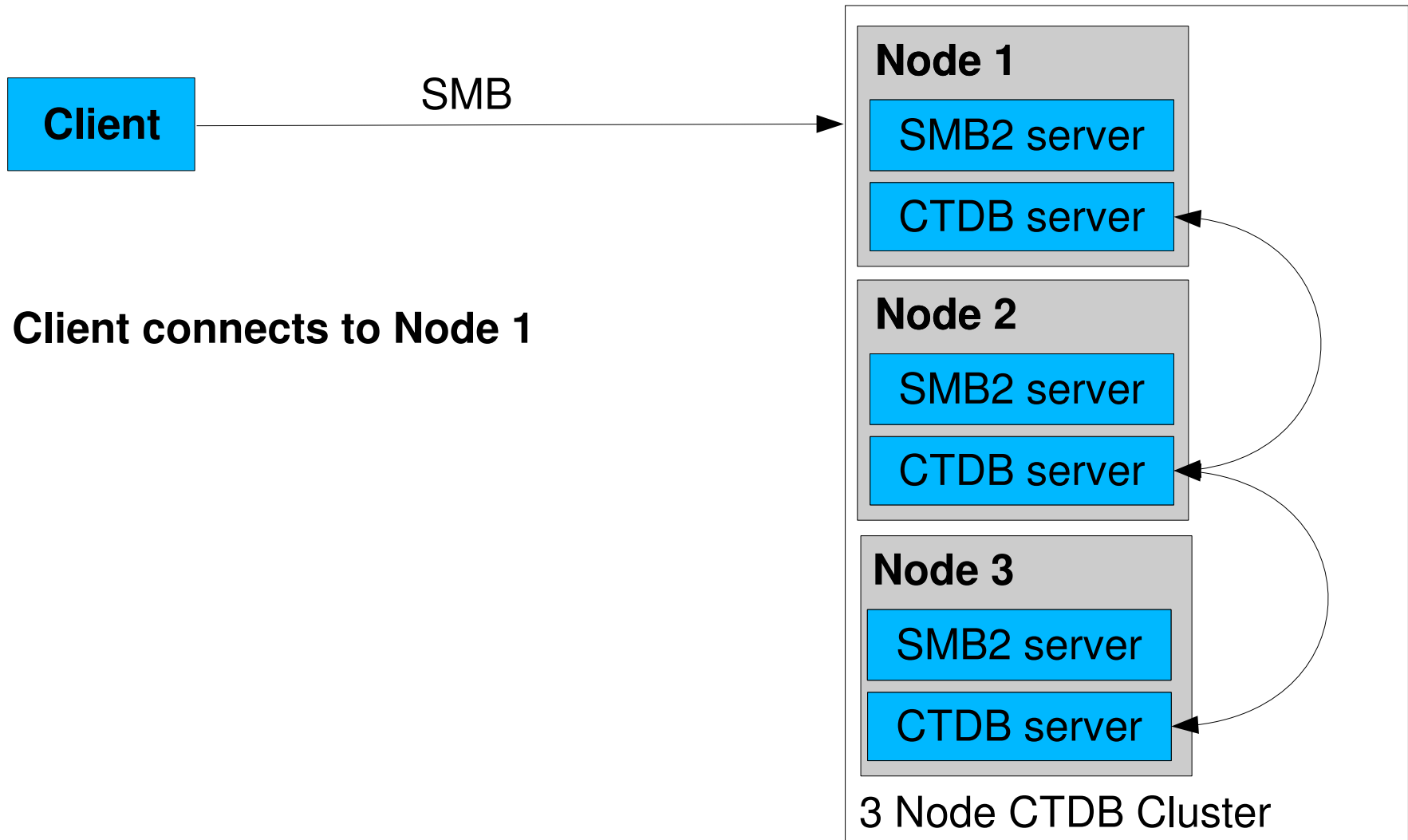
# Failover in SMB1/SMB2 with CTDB

- **CTDB uses Tickle ACKs to speedup recovery**
- **Tickle ACKs:**
  - **are TCP ACK packets with invalid sequence and acknowledge numbers**
  - **cause a TCP client to reestablish a connection with proper sequence numbers, immediately**
  - **were invented/discovered by tridge while working on CTDB**

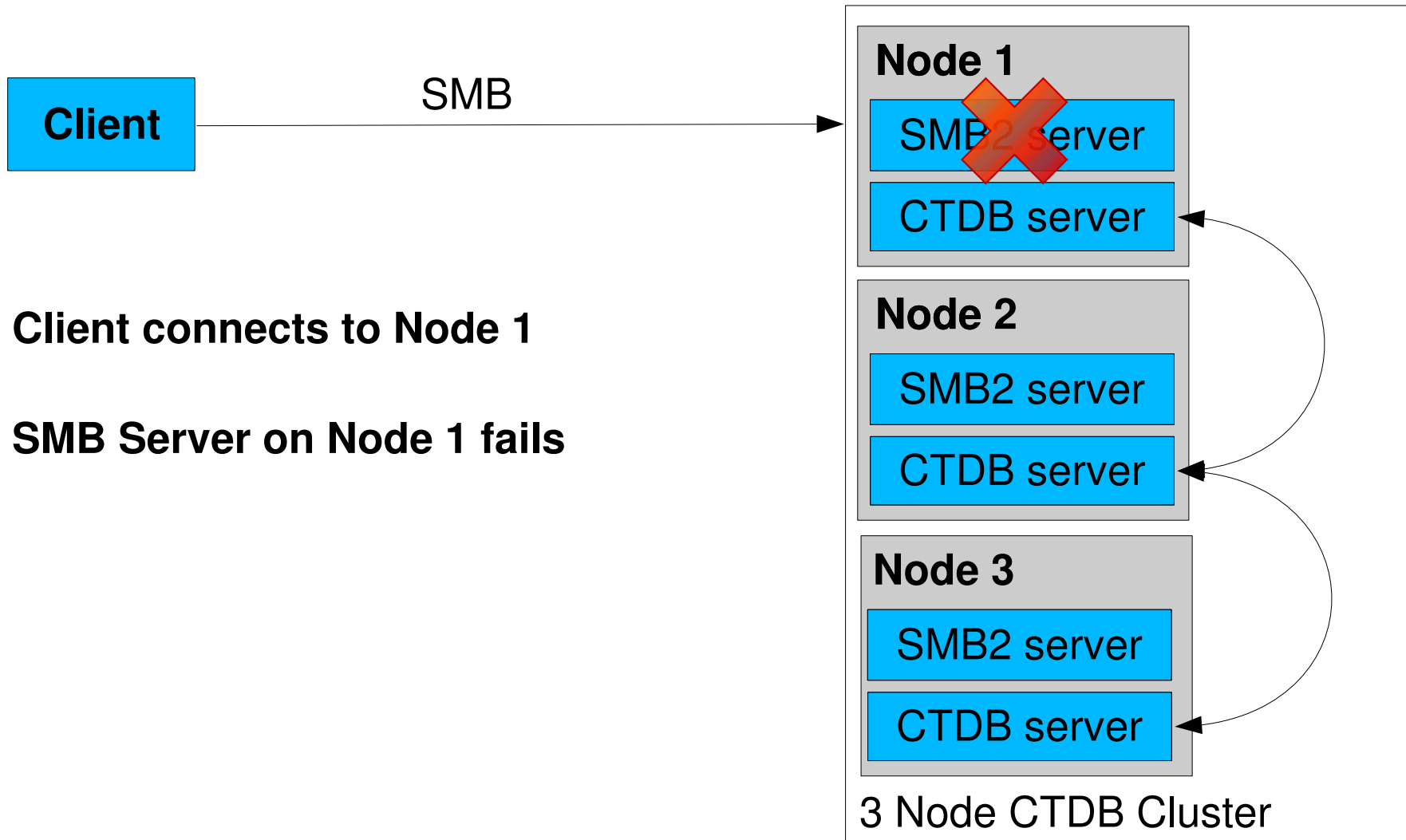
The Pacemaker project also provides a Tickle ACK implementation for use outside of Samba, but that's another presentation.

[Slides](#) and [audio](#) of said presentation. ;) -jarrpa

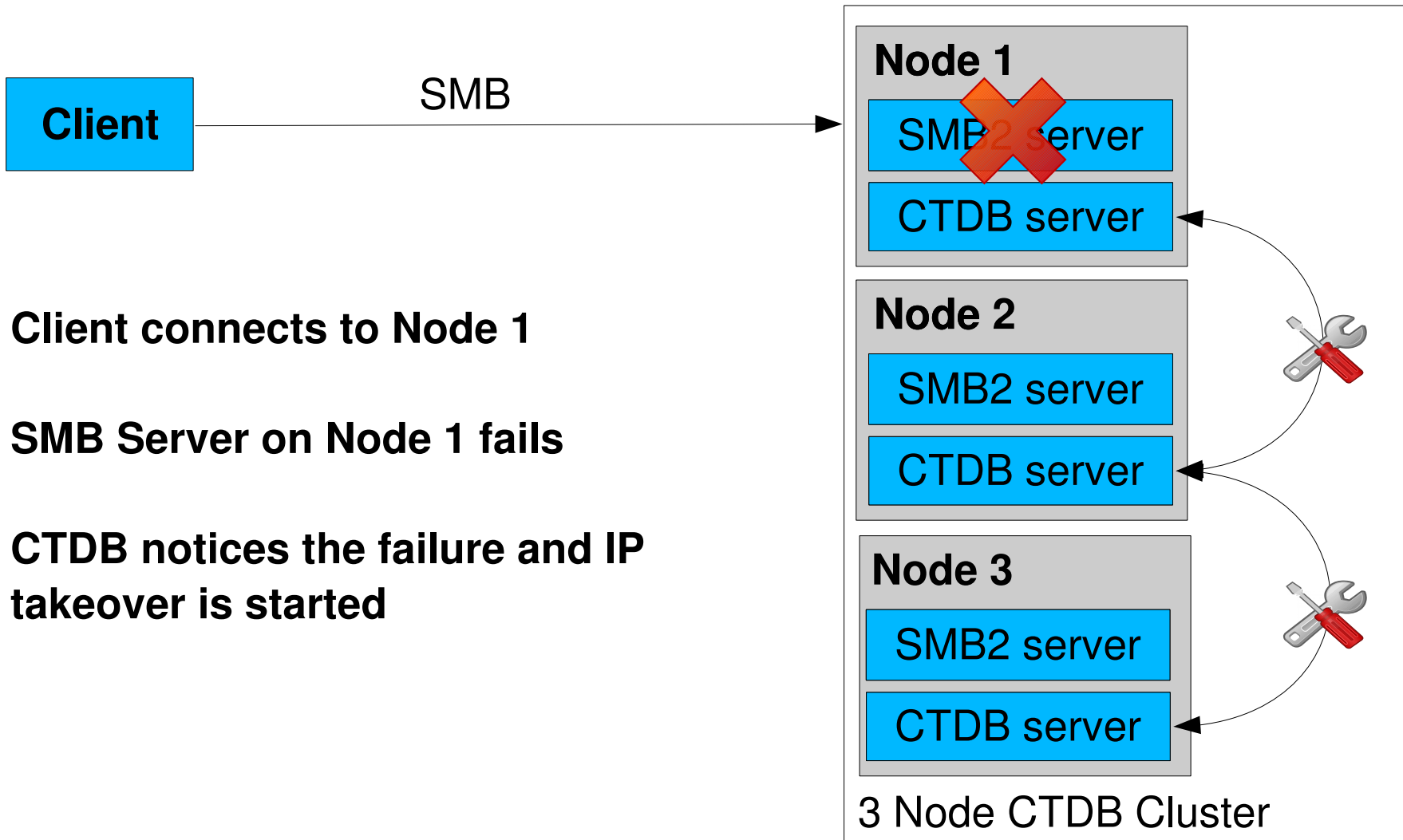
# Failover in SMB1/SMB2 with CTDB



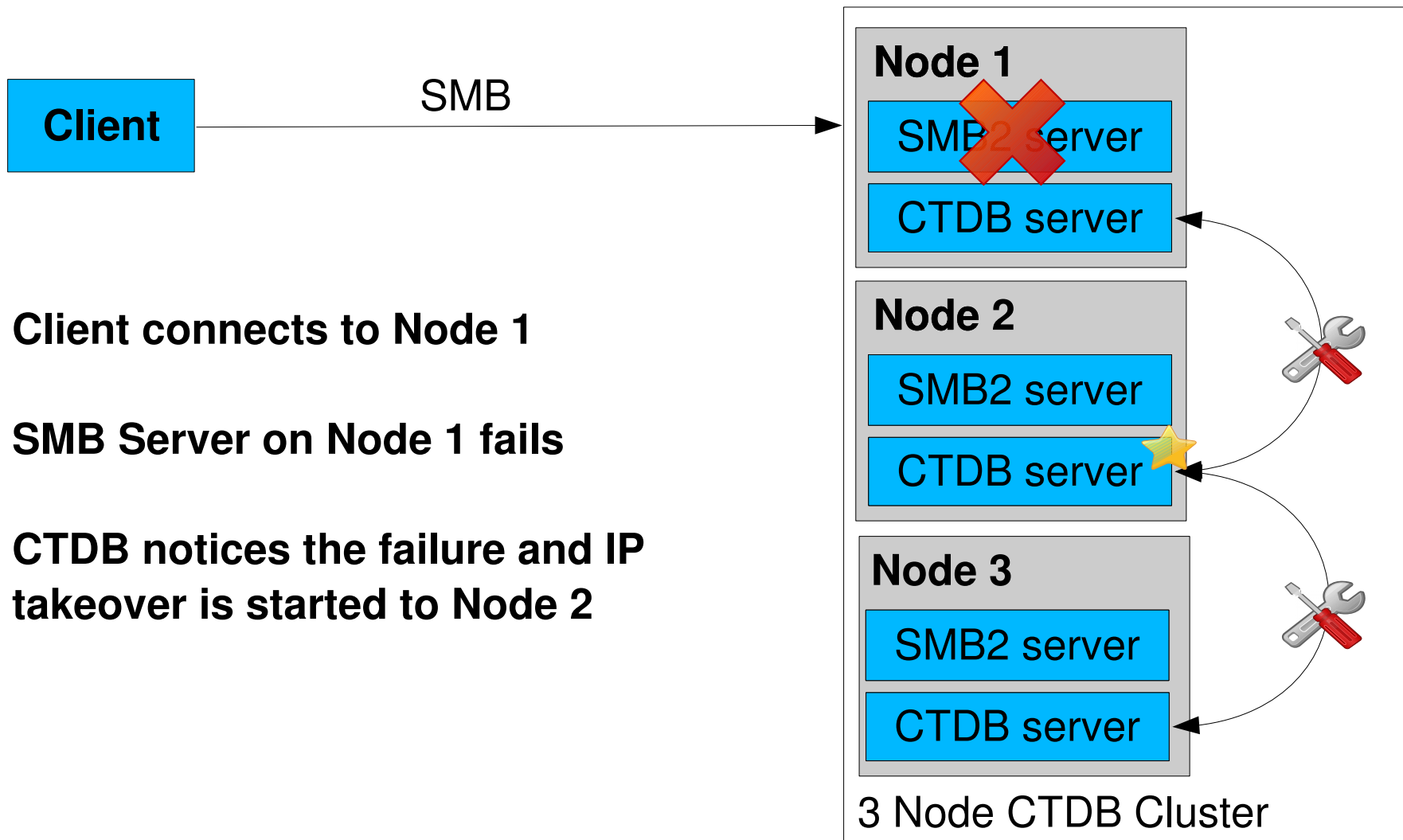
# Failover in SMB1/SMB2 with CTDB



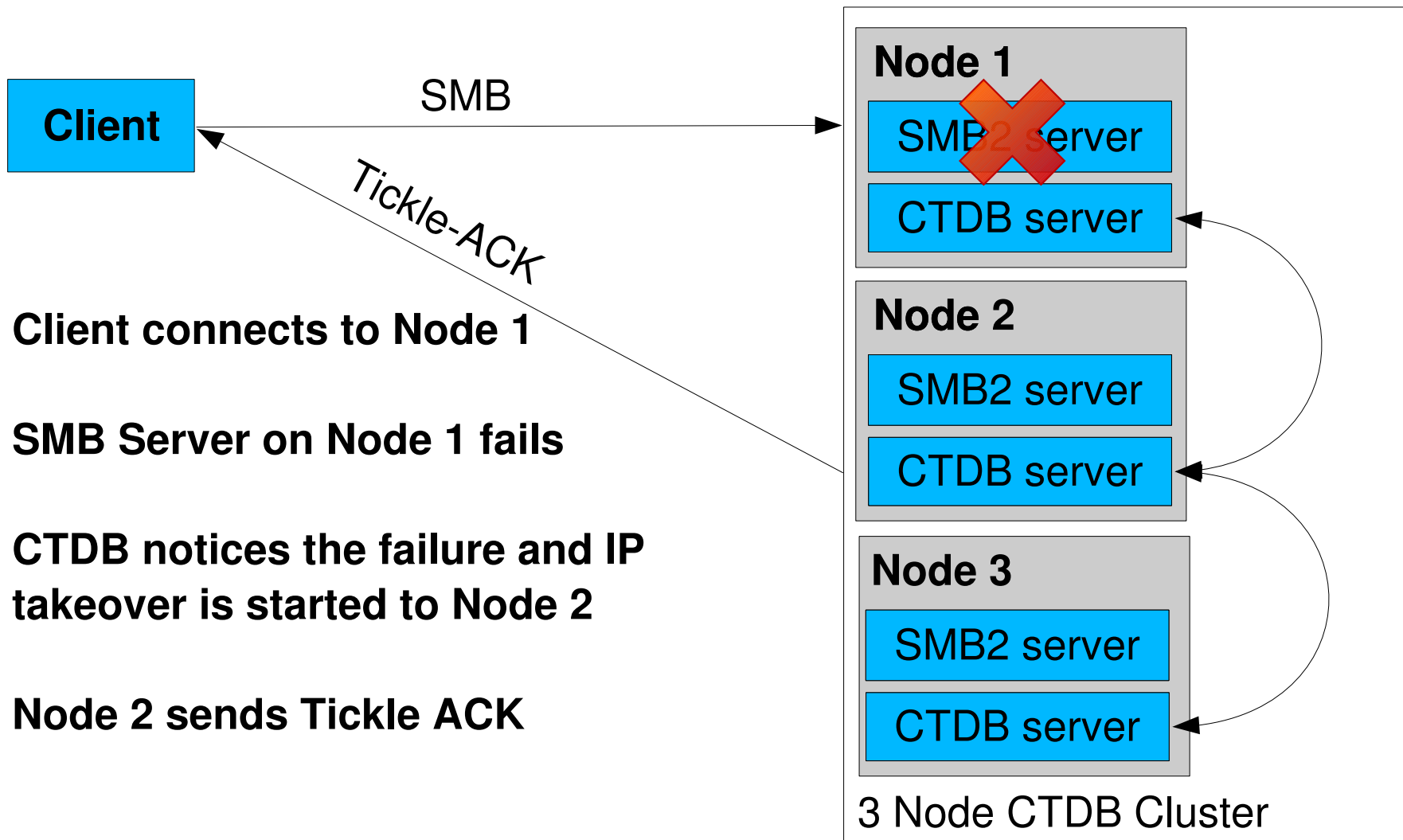
# Failover in SMB1/SMB2 with CTDB



# Failover in SMB1/SMB2 with CTDB

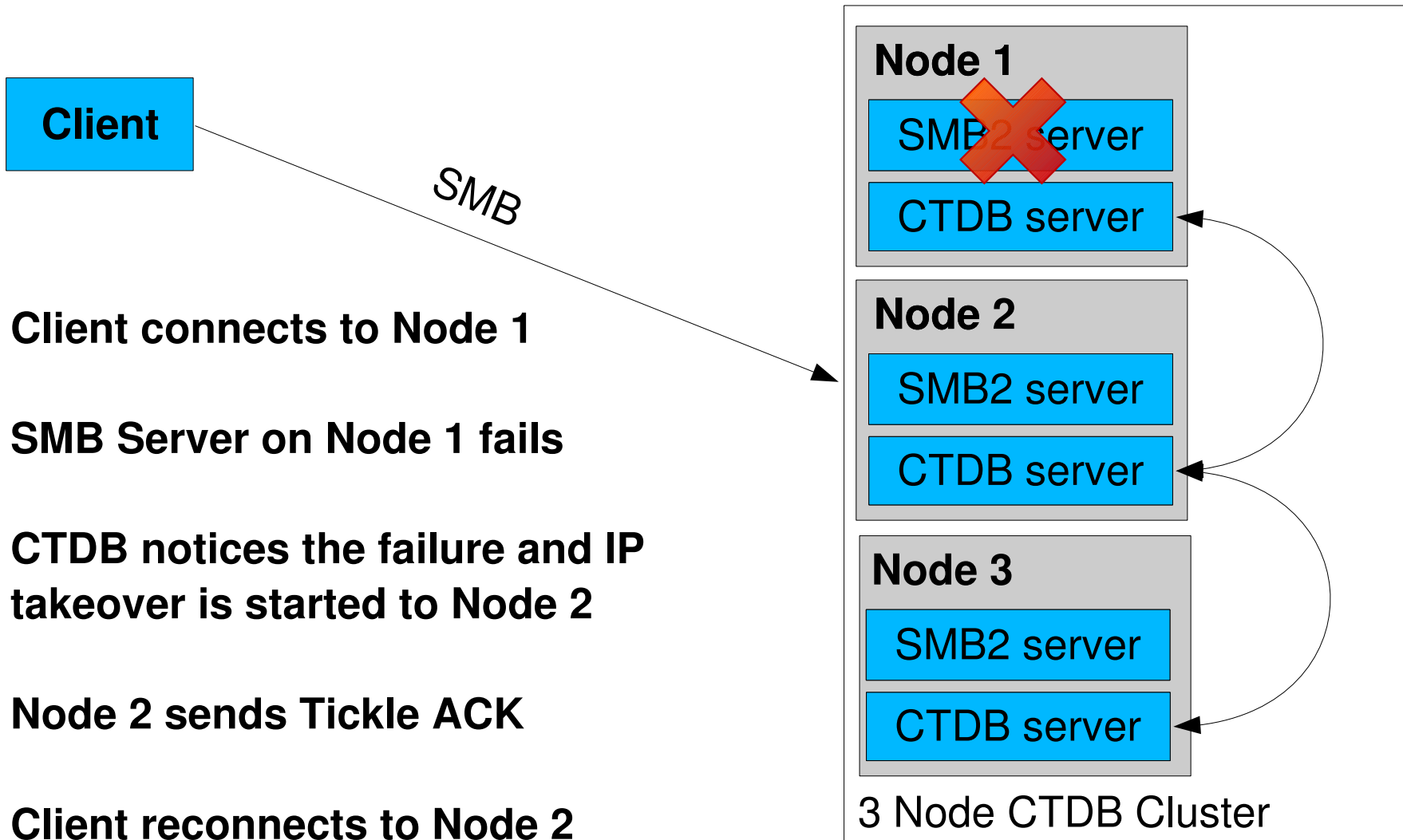


# Failover in SMB1/SMB2 with CTDB





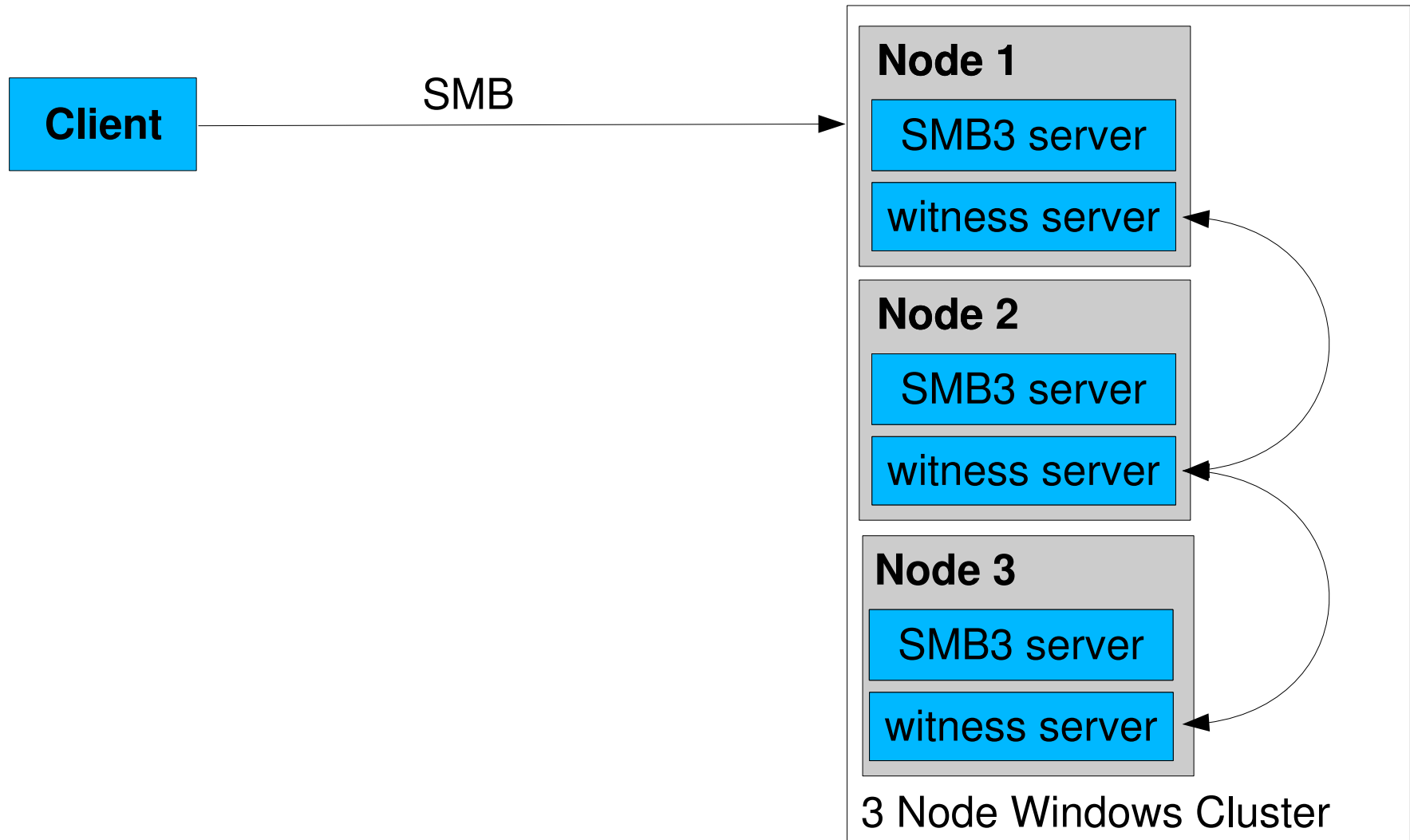
# Failover in SMB1/SMB2 with CTDB



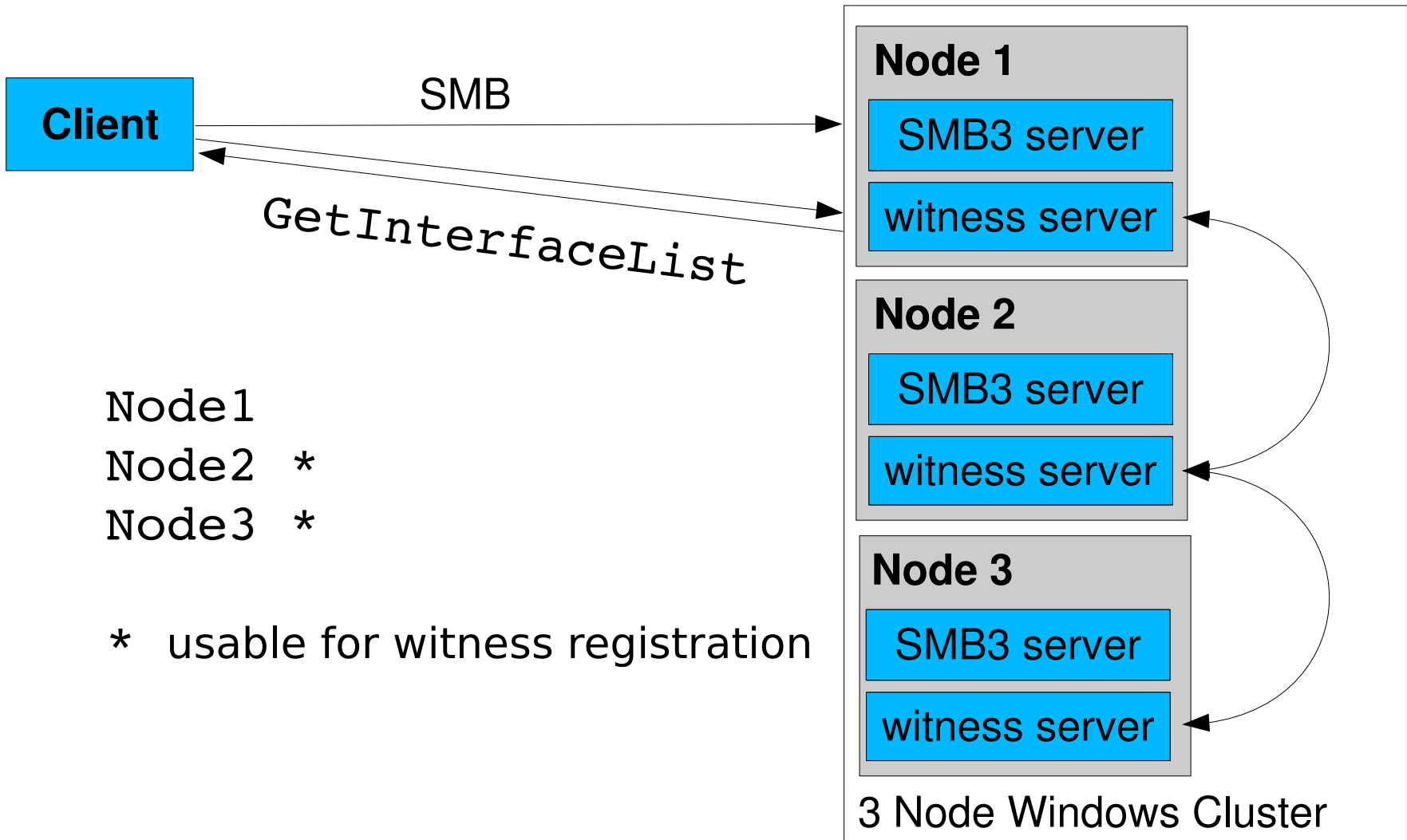
# Failover in SMB3

- **SMB3 achieves transparent failover via several new features:**
  - **Continuous Availability**
  - **Persistent Handles**
  - **Witness**
- **This leads to faster recovery from unplanned node failures**
- **Also allows planned and controlled migration of clients between cluster nodes**

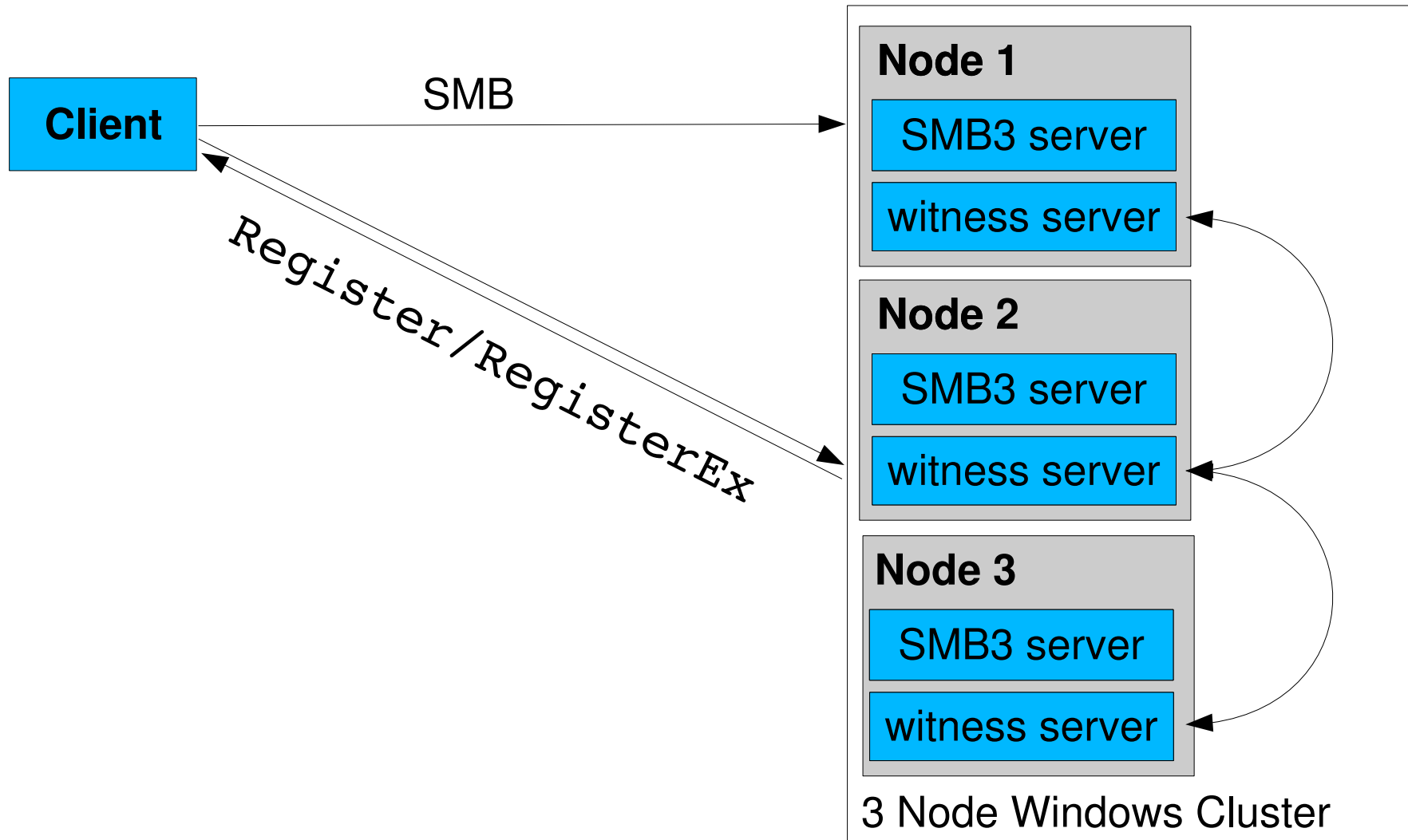
# Failover in SMB3



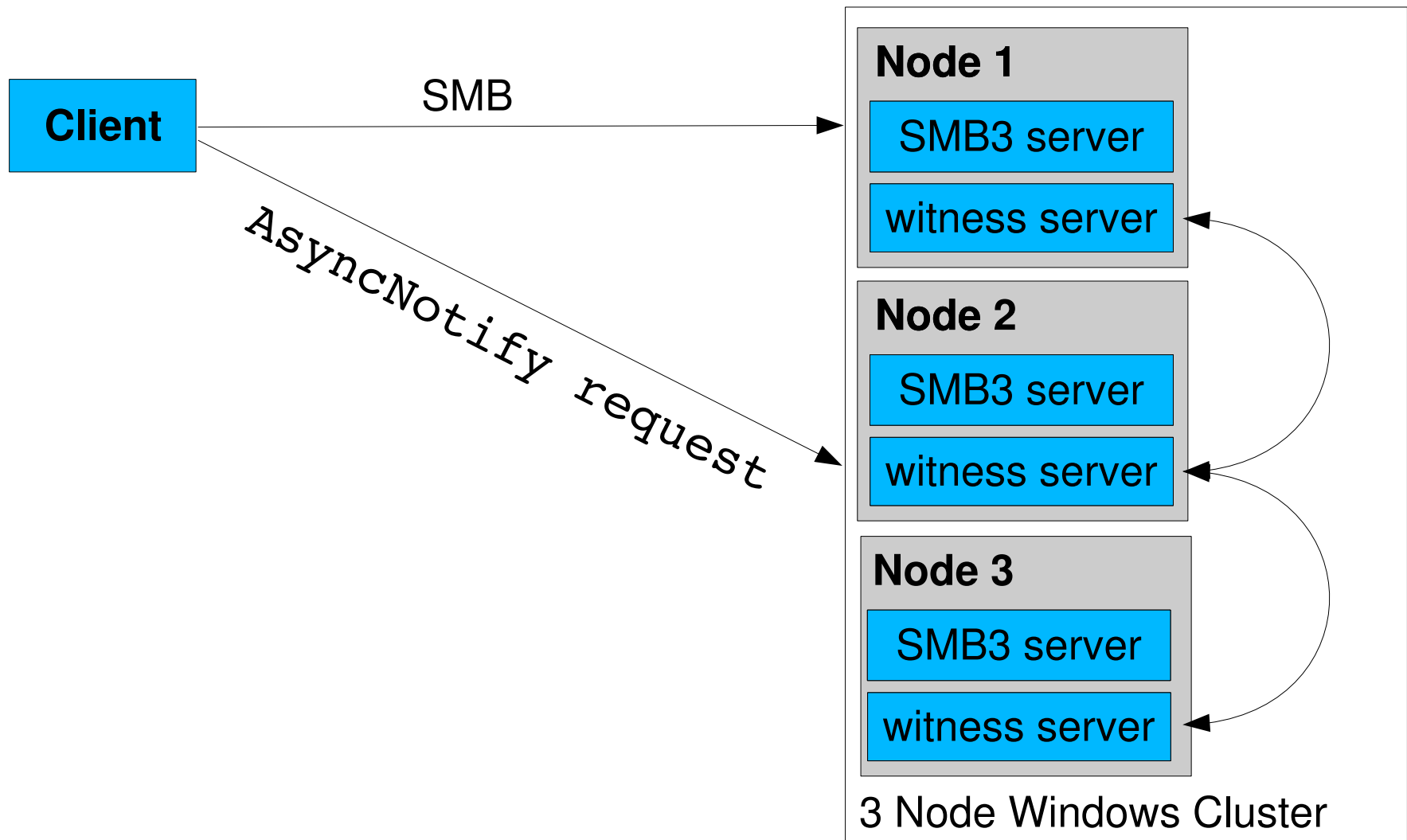
# Failover in SMB3



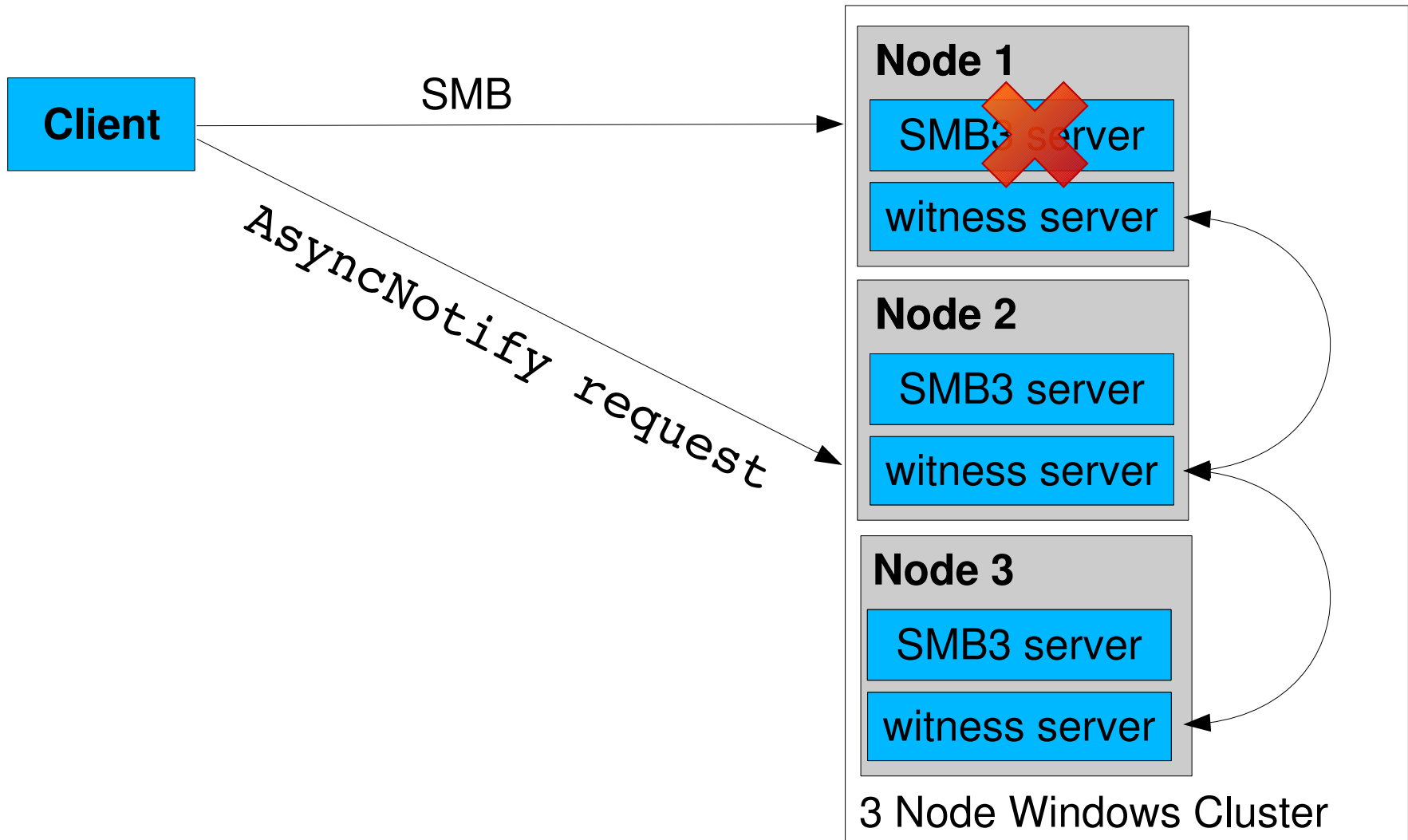
# Failover in SMB3



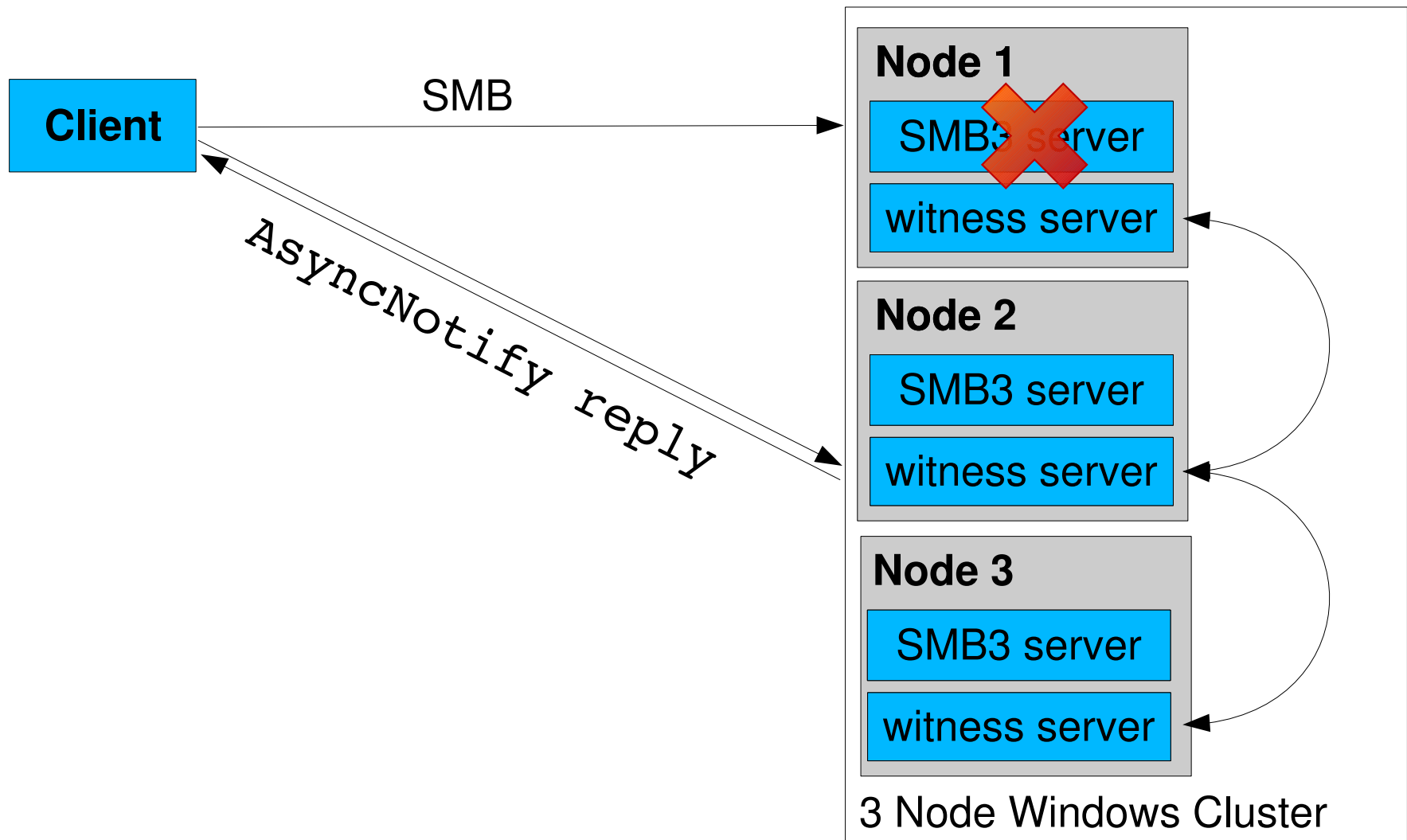
# Failover in SMB3



# Failover in SMB3

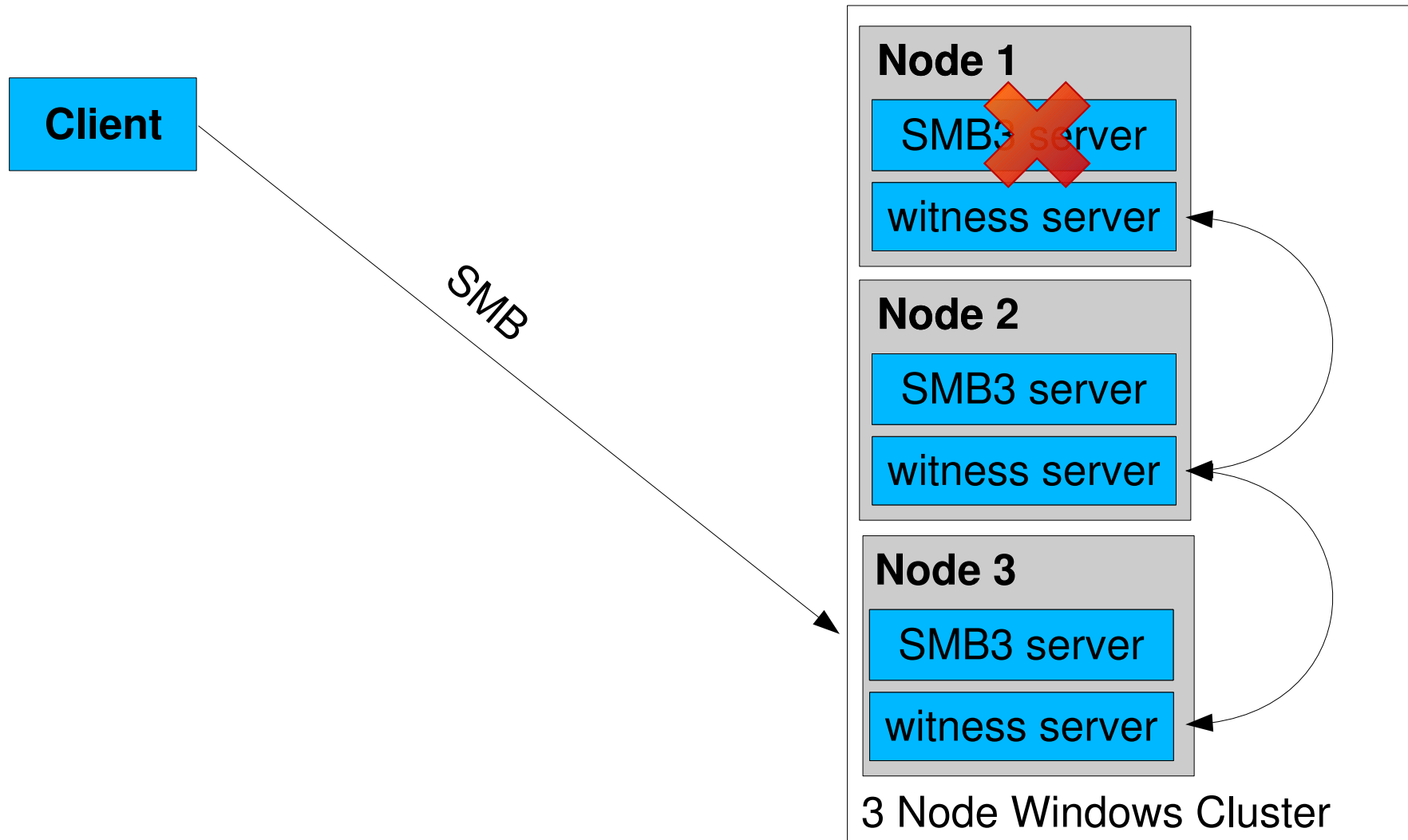


# Failover in SMB3





# Failover in SMB3



## Wait. So why a new protocol?

- **Witness is not strictly about failover, it should be thought of as a notification system**
- **Witness doesn't entirely care who its notifying of what**
- **This allows administrators to programmatically control clients for scenarios such as:**
  - **Load balancing**
  - **Server node maintenance**

# Relationship to SMB3 protocol

- Per share flag enables use of Witness Protocol
- MS-SMB2: “The specified share is present on a server configuration which provides monitoring of the availability of share through the Witness service specified in [MS-SWN]”
- SMB2 TREE\_CONNECT Response Capability Flag:  
SMB2\_SHARE\_CAP\_CLUSTER = 0x00000040
- Witness support seems to be independent from  
SMB2\_SHARE\_CAP\_SCALEOUT and  
SMB2\_SHARE\_CAP\_CONTINUOUS\_AVAILABILITY
- Currently for testing:
  - `smbd:announce CLUSTER = yes`

# The witness interface

- Surprisingly short spec (only 47 pages)
- Version 1, SMB 3.0 (Windows 2012, Windows 8)
- Version 2, SMB 3.02 (Windows 2012 R2, Windows 8.1)
- Only 5 opcodes in the interface:
  - `_witness_GetInterfaceList`
  - `_witness_Register`
  - `_witness_Unregister`
  - `_witness_AsyncNotify`
  - `_witness_RegisterEx` (witness version 2)

# GetInterfaceList

```
DWORD WitnessrGetInterfaceList(  
    [in] handle_t Handle,  
    [out] PWITNESS_INTERFACE_LIST * InterfaceList);
```

- Returns list of network interfaces with IPv4 and/or IPv6 addresses
- Each interface carries information about the interfaces version, state and whether it is a good candidate for witness use

# Witness\_InterfaceInfo

```
interfaces: struct witness_interfaceInfo
  group_name      : 'MTHELENA'
  version         : WITNESS_UNSPECIFIED_VERSION (-1)
  state           : WITNESS_STATE_AVAILABLE (1)
  ipv4            : 192.168.56.108
  ipv6            : ::
  flags           : 0x00000005 (5)
                  1: WITNESS_INFO_IPv4_VALID
                  0: WITNESS_INFO_Ipv6_VALID
                  1: WITNESS_INFO_WITNESS_IF
```

# Register

```
DWORD WitnessrRegister(  
    [in] handle_t Handle,  
    [out] PCONTEXT_HANDLE ppContext,  
    [in] ULONG Version,  
    [in] [string] [unique] LPWSTR NetName,  
    [in] [string] [unique] LPWSTR IPAddress,  
    [in] [string] [unique] LPWSTR ClientComputerName);
```

- Only Witness V1 can be used as version
- Registers client for notify events
- Registration is server-based (NetName) (not share-based)

# UnRegister

```
DWORD WitnessrUnRegister(  
    [in] handle_t Handle,  
    [in] PCONTEXT_HANDLE pContext);
```

- Cleans up client registration



# AsyncNotify

```
DWORD WitnessrAsyncNotify(  
    [in] handle_t Handle,  
    [in] PCONTEXT_HANDLE_SHARED pContext,  
    [out] PRESP_ASYNC_NOTIFY * pResp);
```

- **Asynchronous call**
- **Clients send request and wait, and wait, and wait...**
- **Only in the event of a notification issued by the cluster the client receives a reply**
- **Witness keep-alive mechanism available in Witness v2 (SMB 3.02)**

# AsyncNotify call

- 4 different events are currently defined in the protocol:
- **WITNESS\_NOTIFY\_RESOURCE\_CHANGE**
  - Notify about a resource change state (available, unavailable)
- **WITNESS\_NOTIFY\_CLIENT\_MOVE**
  - Notify a connected client to move to another node
- **WITNESS\_NOTIFY\_SHARE\_MOVE (only v2)**
  - Notify that a share has been moved to another node
- **WITNESS\_NOTIFY\_IP\_CHANGE (only v2)**
  - Notify about an ip address change (online, offline)

# RegisterEx

```
DWORD WitnessrRegisterEx(  
    [in] handle_t Handle,  
    [out] PCONTEXT_HANDLE ppContext,  
    [in] ULONG Version,  
    [in] [string] [unique] LPWSTR NetName,  
    [in] [string] [unique] LPWSTR ShareName,  
    [in] [string] [unique] LPWSTR IpAddress,  
    [in] [string] [unique] LPWSTR ClientComputerName,  
    [in] ULONG Flags,  
    [in] ULONG KeepAliveTimeout);
```

- Available with Windows 2012 R2 (Witness v2)
- Witness keepalive as client can define KeepAliveTimeout
- Server returns with ERROR\_TIMEOUT after KeepAliveTimeout has expired (Windows 8.1 default 120 seconds)

# RegisterEx

- **Optional ShareName allows share notify instead of server notify**
- **Allows for use of asymmetric storage (SMB 3.02)**
- **Flags field allows tracking of IP notifications**

# witness testing

- **rpcclient witness command set**
- **smbtorture local.ndr.witness**
  - **Just tests correctness of the NDR marshalling/unmarshalling**
- **smbtorture rpc.witness**
  - **Test correctness of the DCE/RPC calls**
- **Fundamental problem: how to test a cluster ? How to test resource changes? How to test node failures ?**
- **Windows Failover Cluster Manager does resource changes with yet another DCE/RPC protocol**

# Sidetrack: clusapi

- **Microsoft Cluster Management API**
  - > 200 opcodes
  - > 600 pages protocol spec
  - Used by Microsoft Failover Cluster Manager
- purely DCE/RPC based interface (over ncacn\_ip\_tcp[seal])
- Samba now has IDL (for v3 of that protocol) and a torture test suite
- MS-CRMP  
Failover Cluster: Management API (ClusAPI) Protocol
- Some ideas to use this protocol as front-end for remote CTDB management

# Sidetrack: clusapi

- **Basic CLUSAPI v3 implementation in Samba**
- **“Failover Cluster Manager” on Windows insists on contacting DCOM interfaces which Samba currently does not support**
- **Spme Cluster Power Shell commandlets already work against Samba**
- **Current WIP branch:**  
<https://git.samba.org/?p=gd/samba/.git;a=shortlog;h=refs/heads/master-clusapi>

# DCE/RPC requirements

- endpointmapper with ncacn\_ip\_tcp support
- DCE/RPC sign & seal (SPNEGO,KRB5,NTLMSSP)
- asynchronous DCE/RPC server
  - Currently two unfinished implementations:
    - David Disseldorp <[ddiss@samba.org](mailto:ddiss@samba.org)>
    - Stefan Metzmacher <[metze@samba.org](mailto:metze@samba.org)>
  - (also needed for MS-PAR and possibly other protocols)
- mgmt service (Remote DCE/RPC service management)
  - Two implementations available, none is published yet.
  - `mgmt_inq_princ_name()` for different node principals



# witnessd server

- **Standalone binary, using new infrastructure invented for spoolssd**
- **Independent binary so any Samba server problem does not interfere with witness messaging**
- **Needs to register for at least 4 notification events (messaging)**
- **Configuration and possibly Server State store**
- **Very close integration with CTDB:**
  - **CTDB maintains all available cluster state information**
  - **CTDB already has mechanisms to communicate failures between the nodes**
  - **CTDB could easily reuse tickle-ack hooks for witness notifications**

# Roadmap: Witness support in Samba

- **Early PoC implementation by Gregor Beck and Stefan Metzmacher from 2012**
- **Wireshark dissector for witness protocol (not upstream yet)**
- **Full IDL and torture tests in Samba Git repository upstream**
- **Witness Service is on Samba Roadmap as a funded project**
  - **Goal: Samba 4.4/4.5 should have a full witness implementation**
- **Currently resolving some infrastructure requirements**

# Roadmap: Witness client interface

- **Frontend for management tasks of witness server:**
- **listing of active, connected clients  
(shared state stored in distributed database)**
- **Tool to manually move Clients to other nodes  
(similar to Move-SmbWitnessClient PowerShell cmdlet)**
- **Tool to move share to other node**
- **Currently implemented as part of the smbcontrol management and messaging tool**
  - **“smbcontrol witnessd witnessnotify”**
  - **subcommands: change, move, sharemove, ipnotify**

# Roadmap: Integration w/external projects

- **Several existing SMB clients would benefit from supporting Witness, including:**
  - CIFS Kernel module
  - smbclient
  - libsmbclient
- **Alternatives to CTDB could be introduced for the purposes of tracking the state of resources in the cluster (e.g. Pacemaker).**
- **Possibly implement a stand-alone Witness client service to monitor/witness non-SMB connections.**

# Further reading

- **Microsoft Protocol Documentation:**
  - **MS-SWN: Service Witness Protocol**
  - **MS-SMB2: Server Message Block (SMB) Protocol Versions 2 and 3**
  - **MS-CMRP: Failover Cluster Management Protocol**
- **SMB 2.x and SMB 3.0 Timeouts in Windows**  
<http://blogs.msdn.com/b/openspecification/archive/2013/03/27/smb-2-x-and-smb-3-0-timeouts-in-windows.aspx>
- **Samba Wiki**  
[https://wiki.samba.org/index.php/Samba3/SMB2#Witness\\_Notification\\_Protocol](https://wiki.samba.org/index.php/Samba3/SMB2#Witness_Notification_Protocol)

# Questions and answers

- Mail [gd@samba.org](mailto:gd@samba.org) or [jarrpa@samba.org](mailto:jarrpa@samba.org)
- #samba-technical on irc.freenode.net
- Latest stable WIP:  
<https://git.samba.org/?p=gd/samba/.git;a=shortlog;h=refs/heads/master-witness-ok>

# Thank you for your attention!

[www.redhat.com](http://www.redhat.com)  
[www.samba.org](http://www.samba.org)

<[gd@samba.org](mailto:gd@samba.org)>

<[jarrpa@samba.org](mailto:jarrpa@samba.org)>