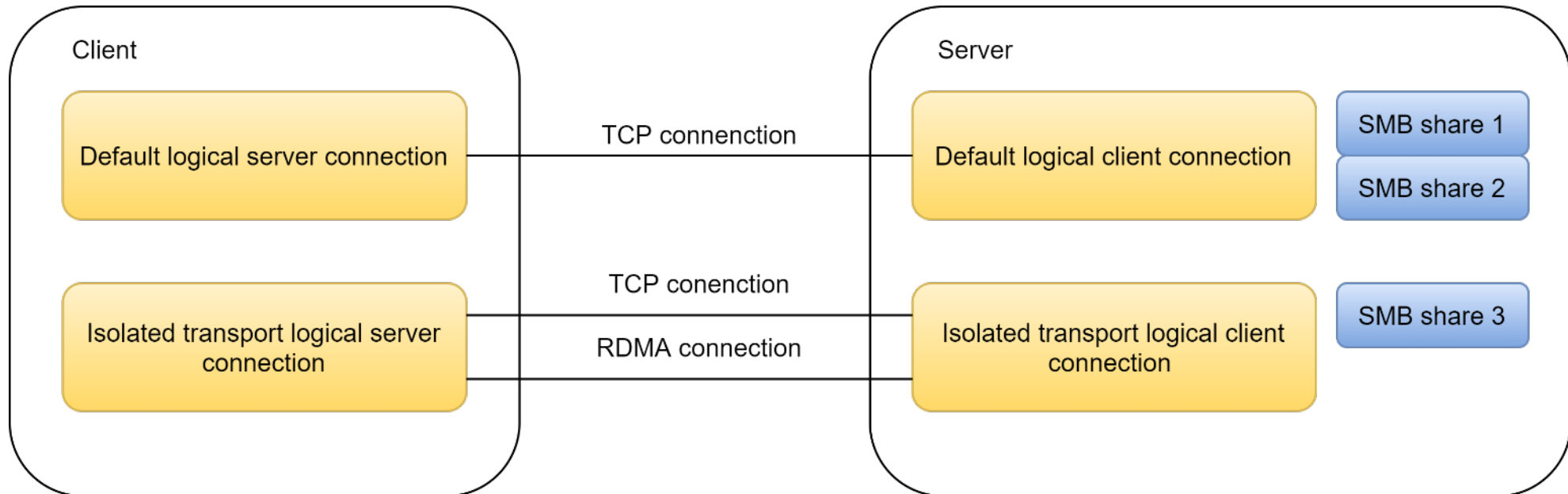# Agenda

- **Isolated Transport**
  - 2 Logical Connections To The Same Server
- **RDMA Transforms**
  - Encryption
  - Signing
  - Compression
- **New Management Options**
  - Encryption Cipher Orders [Ni]
  - Compression Opt-in [Fe]
  - Compressibility Sampling [Ni]
- **Improvements & Bug Fixes**
  - Remote Desktop Services (RDS) AutoDisconnectTimeout
  - Message ID On Cancel Request
  - Live-lock With Memory-Mapped Section Creation In SMB Redirector

# Isolated Transports

STORAGE DEVELOPER CONFERENCE

SDC 21

# Purpose and Design

- Optionally create a secondary logical connection to the same server to isolate transports used (e.g. RDMA vs. no RDMA).

# RDMA Transforms

# Protocol Changes - Negotiation

## 2.2.3.1.6 SMB2_RDMA_TRANSFORM_CAPABILITIES

The SMB2_RDMA_TRANSFORM_CAPABILITIES context is specified in an SMB2 NEGOTIATE request by the client to indicate the transforms supported when data is sent over RDMA. The format of the data in the **Data** field of this SMB2_NEGOTIATE_CONTEXT is as follows:

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 1 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 2 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 3 0 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| TransformCount | | | | | | | | | | | | | | | | Reserved1 | | | | | | | | | | | | | | | |
| Reserved2 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| RDMATransformIds (variable) | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| ... | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

**TransformCount (2 bytes):** The number of elements in **RDMATransformIds** array. This value MUST be greater than 0.
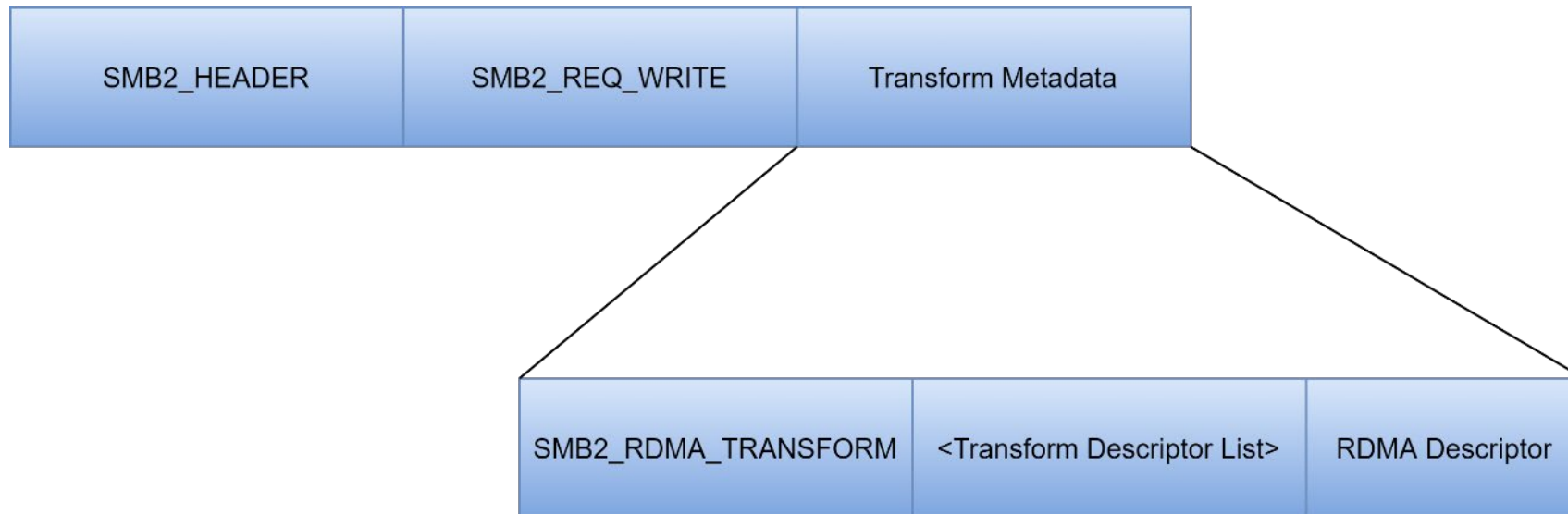
**Reserved1 (2 bytes):** This field MUST NOT be used and MUST be reserved. The sender MUST set this to 0, and the receiver MUST ignore it on receipt.

**Reserved2 (4 bytes):** This field MUST NOT be used and MUST be reserved. The sender MUST set this to 0, and the receiver MUST ignore it on receipt.

**RDMATransformIds (variable):** An array of 16-bit integer IDs specifying the supported RDMA transforms. The following IDs are defined.

STORAGE DEVELOPER CONFERENCE
SDC 21

# Protocol Changes - Packet Diagram

- E.g. SMB2 Write over RDMA

| SMB2_HEADER | SMB2_REQ_WRITE | Transform Metadata |
|---|---|---|

| SMB2_RDMA_TRANSFORM | <Transform Descriptor List> | RDMA Descriptor |
|---|---|---|

# Protocol Changes – Transform Metadata

## 2.2.43 SMB2_RDMA_TRANSFORM

The SMB2_RDMA_TRANSFORM is used by the client or server to send/receive transformed RDMA payload in READ/WRITE operations. The SMB2_RDMA_TRANSFORM is optional and only valid for the SMB 3.1.1 dialect when connection supports RDMA transform.<82>

| 0 1 2 3 4 5 6 7 8 9 | 1 0 1 2 3 4 5 | 6 7 8 9 | 2 0 1 2 3 4 5 6 7 8 9 | 3 0 1 |
|---|---|---|---|---|

| RdmaDescriptorOffset | RdmaDescriptorLength |
|---|---|
| Channel | |
| TransformCount | Reserved1 |
| Reserved2 | |

**RdmaDescriptorOffset (2 bytes):** This field contains the offset, in bytes, from the beginning of this structure to the RDMA descriptors as specified by the **Channel** field.

**RdmaDescriptorLength (2 bytes):** This field contains the length, in bytes, of the RDMA descriptors as specified by the **Channel** field.

**Channel (4 bytes):** This field MUST contain exactly one of the following values:

| Value | Meaning |
|---|---|
| SMB2_CHANNEL_NONE<br>0x00000000 | No channel information is present. The **RdmaDescriptorOffset** and **RdmaDescriptorLength** fields MUST be set to zero by the sender and MUST be ignored by the receiver. |
| SMB2_CHANNEL_RDMA_V1<br>0x00000001 | One or more SMB_DIRECT_BUFFER_DESCRIPTOR_V1 structures as specified in [MS-SMBD] section 2.2.3.1 are present in the channel information specified by the **RdmaDescriptorOffset** and **RdmaDescriptorLength** fields. |
| SMB2_CHANNEL_RDMA_V1_INVALIDATE<br>0x00000002 | One or more SMB_DIRECT_BUFFER_DESCRIPTOR_V1 structures as specified in [MS-SMBD] section 2.2.3.1 are present in the channel information specified by the **RdmaDescriptorOffset** and **RdmaDescriptorLength** fields. The server is requested to perform remote invalidation when responding to the request as specified in [MS-SMBD] section 3.1.4.2. |

**TransformCount (2 bytes):** This field specifies the number of transforms present after this structure. This value MUST be greater than 0.

## 2.2.43.1 SMB2_RDMA_CRYPTO_TRANSFORM

The SMB2_RDMA_CRYPTO_TRANSFORM is used by the client or server to send/receive encrypted or signed RDMA payload in READ/WRITE operations. The SMB2_RDMA_CRYPTO_TRANSFORM is optional and only valid for the SMB 3.1.1 dialect.<83>

| 0 1 2 3 4 5 6 7 8 9 | 1 0 1 2 3 4 5 | 6 7 8 9 | 2 0 1 2 3 4 5 6 7 8 9 | 3 0 1 |
|---|---|---|---|---|

| TransformType | SignatureLength |
|---|---|
| NonceLength | Reserved |
| Signature (variable) | |
| ... | |
| Nonce (variable) | |
| ... | |
| Padding (variable) | |

**TransformType (2 bytes):** This field MUST be set to one of the following values.

| Value | Meaning |
|---|---|
| SMB2_RDMA_TRANSFORM_TYPE_ENCRYPTION<br>0x0001 | RDMA transform of type encryption is present and the payload is encrypted. |
| SMB2_RDMA_TRANSFORM_TYPE_SIGNING<br>0x0002 | RDMA transform of type signing is present and the payload is signed. |

**SignatureLength (2 bytes):** The length, in bytes, of **Signature** field.

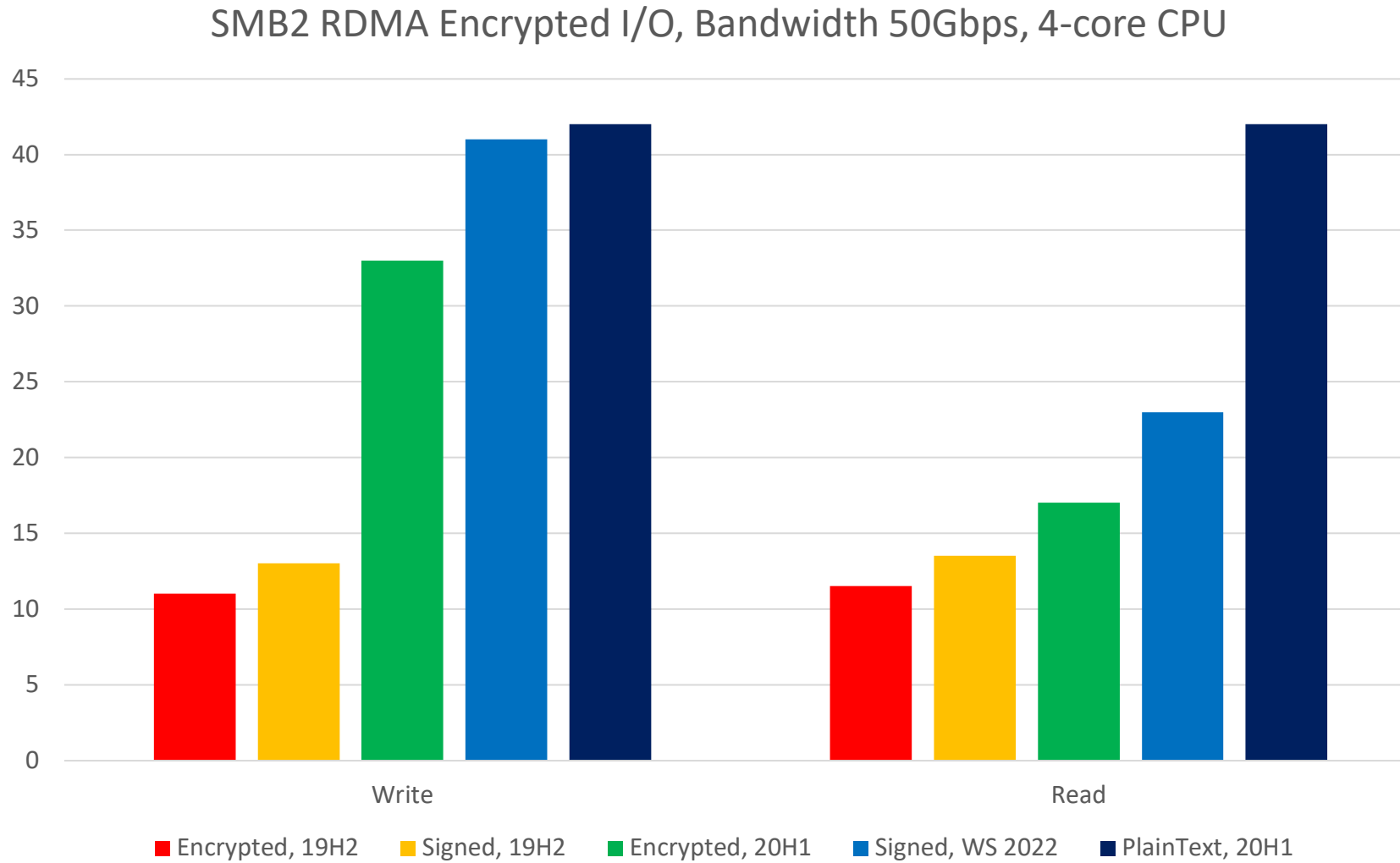**NonceLength (2 bytes):** The length, in bytes, of **Nonce** field.

**Reserved (2 bytes):** This field MUST NOT be used and MUST be reserved. The sender MUST set this to zero, and the receiver MUST ignore it on receipt.

**Signature (variable):** The signature of the encrypted/signed data generated using **Session.EncryptionKey**. The length of this field MUST be less than or equal to 16 bytes.

**Nonce (variable):** An implementation-specific value assigned for encrypted/signed data. This MUST NOT be reused for an SMB2 message within a session.

**Padding (variable):** This optional field is present after **Nonce** field so the channel information, if any, after this structure starts at the first 8-byte aligned offset. The sender MUST set this to zero, and the receiver MUST ignore it on receipt.

STORAGE DEVELOPER CONFERENCE
SDC 21

# Performance



SMB2 RDMA Encrypted I/O, Bandwidth 50Gbps, 4-core CPU

Legend: Encrypted, 19H2 · Signed, 19H2 · Encrypted, 20H1 · Signed, WS 2022 · PlainText, 20H1

# New Management Options

# Encryption Ciphers

- Effects & Expectation Of New Settings
  - `Set-Smb(Client|Server)Configuration –EncryptionCiphers <ciphers>`
  - `Get-Smb(Client|Server)Configuration`
- Current Supported List Of Encryption Ciphers
  - AES_128_GCM
  - AES_256_GCM
  - AES_128_CCM
  - AES_256_CCM
- Race Conditions Considered
  - Negotiate Scenario
  - Setting Ciphers With Multiple Client / Server Instances

# Compression Opt-In

- SMB Server
  - New-SmbShare –CompressData <bool>
  - Set-SmbServerConfiguration –DisableCompression <bool>

- SMB Client
  - New-SmbMapping –CompressData <bool>
  - net use /COMPRESS <bool>

STORAGE DEVELOPER CONFERENCE

SDC 21

# Compressibility Sampling

- Default behavior changed from always sample 512MB to always compress. Feature needs to be opted-in from registry key.
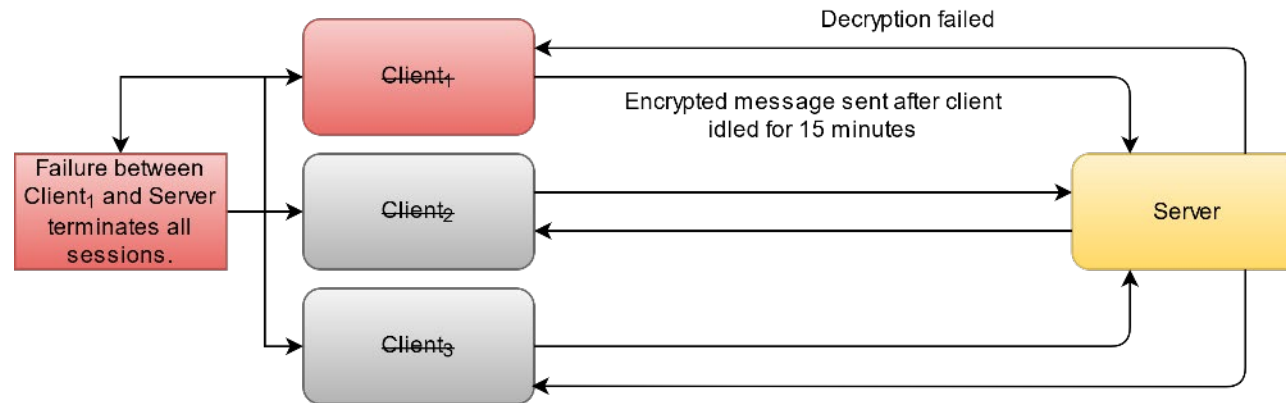
# Improvements & Bug Fixes

STORAGE DEVELOPER CONFERENCE

SDC 21

# RDS AutoDisconnectTimeout

- Background, Issue Context, & Brief Summary
- Solutions Discussed
- Proposed Solution

# RDS AutoDisconnectTimeout

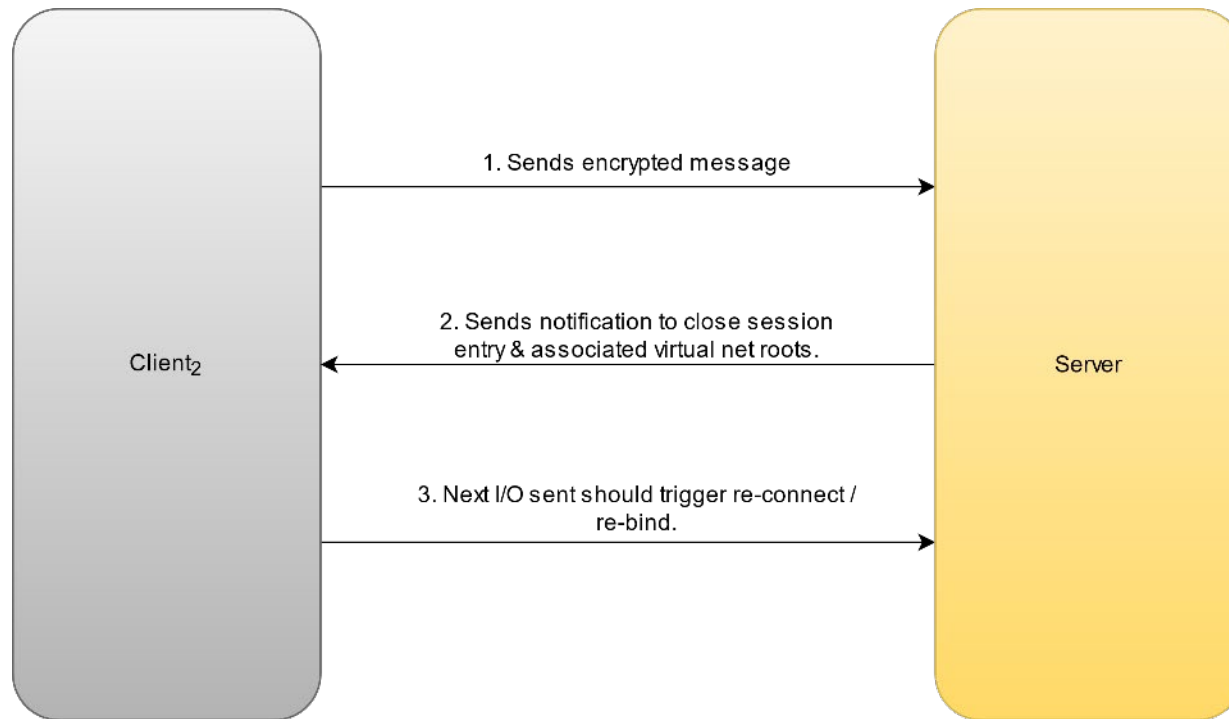- Background, Issue Context, & Brief Summary

# RDS AutoDisconnectTimeout

- ## Solutions Discussed
  - Retaining encryption keys after session destruction
  - No idle session disconnects when connection shared between multiple sessions
  - Including MessageID within the SMB_TRANSFORM_HEADER
  - Notification from server to client like oplock break notification to disconnect the server entry on the client

# RDS AutoDisconnectTimeout

- Proposed Improvement
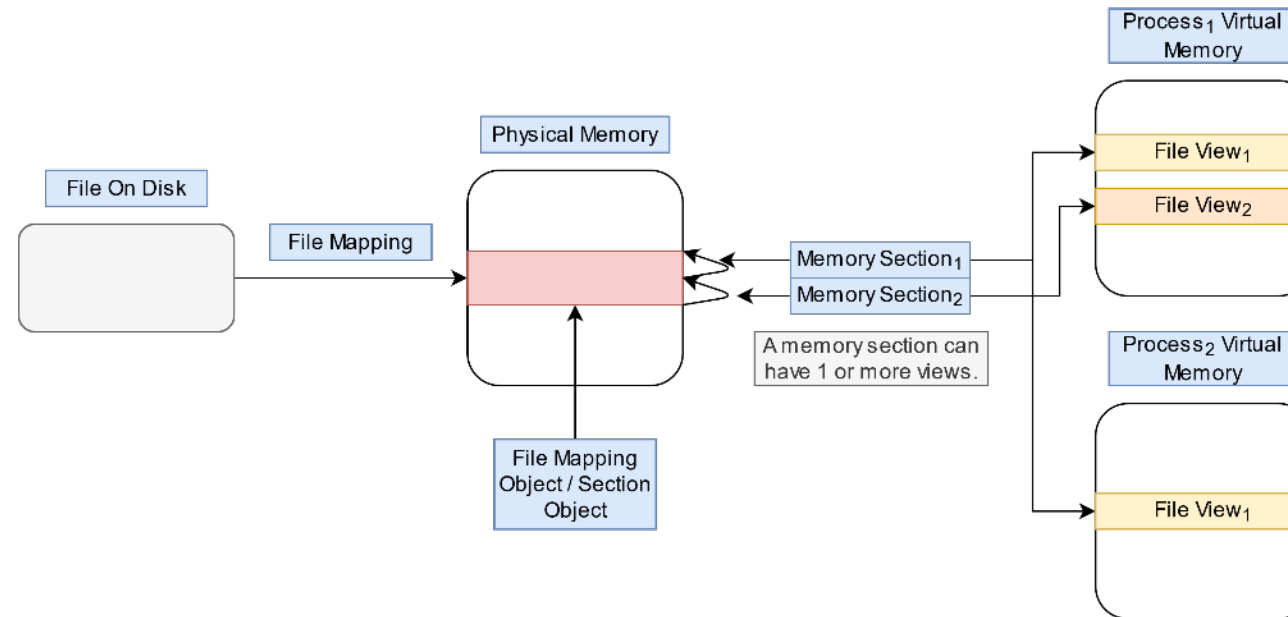
# Message ID on cancel request

- When using AES-GMAC we require a nonce to calculate signature, and message ID is included in the nonce. When sending a cancel request to an async packet, we used to send 0 as the message ID, which violated the nonce's requirement that it should never repeat itself.

# Live-lock With Memory-Mapped Section Creation In SMB Redirector

- Issue Context & Brief Summary
- Background Of Memory-Mapped Files
- Detailed Overview Of Issue
- Solution

# Live-lock



**File Mapping**

Physical Memory

File On Disk

File Mapping

Process$_1$ Virtual Memory

File View$_1$

File View$_2$

Memory Section$_1$
Memory Section$_2$

A memory section can have 1 or more views.

File Mapping Object / Section Object

Process$_2$ Virtual Memory

File View$_1$

**Useful Concepts**

| File Size: Known as EOF where the size of data is in bytes. | Valid Data Length: VDL is less than or equal to EOF; this is the data in the file that's written to | Allocation Size: Size allocated for files on disk; even multiple of cluster size. |
|---|---|---|

STORAGE DEVELOPER CONFERENCE
SDC 21

# Live-lock

- Overview of Issue
- $Client_1$ opens file with shared read & write permissions
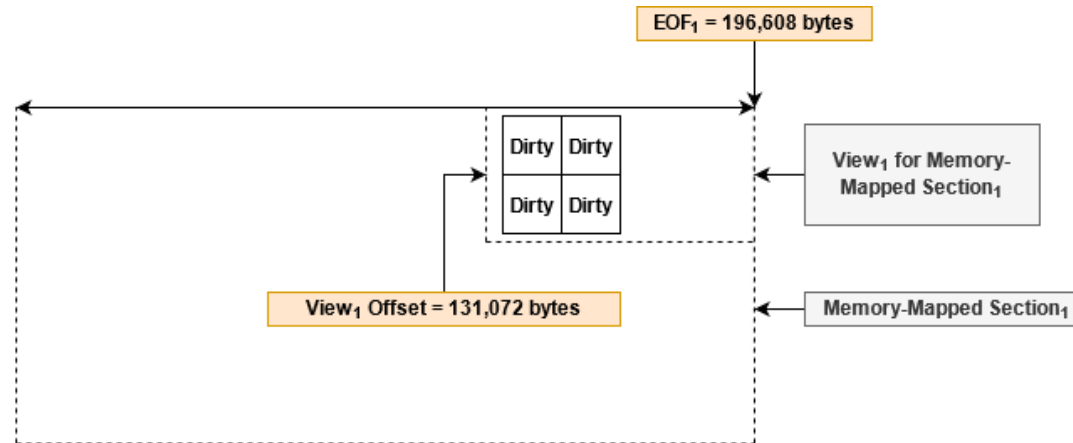  - Obtains RWH oplock
  - Assume empty file to work with

$EOF_1 = 0$ bytes

# Live-lock

- Overview of Issue
- $Client_1$ sets end of file ($EOF_1$) to e.g., 196,608 bytes.
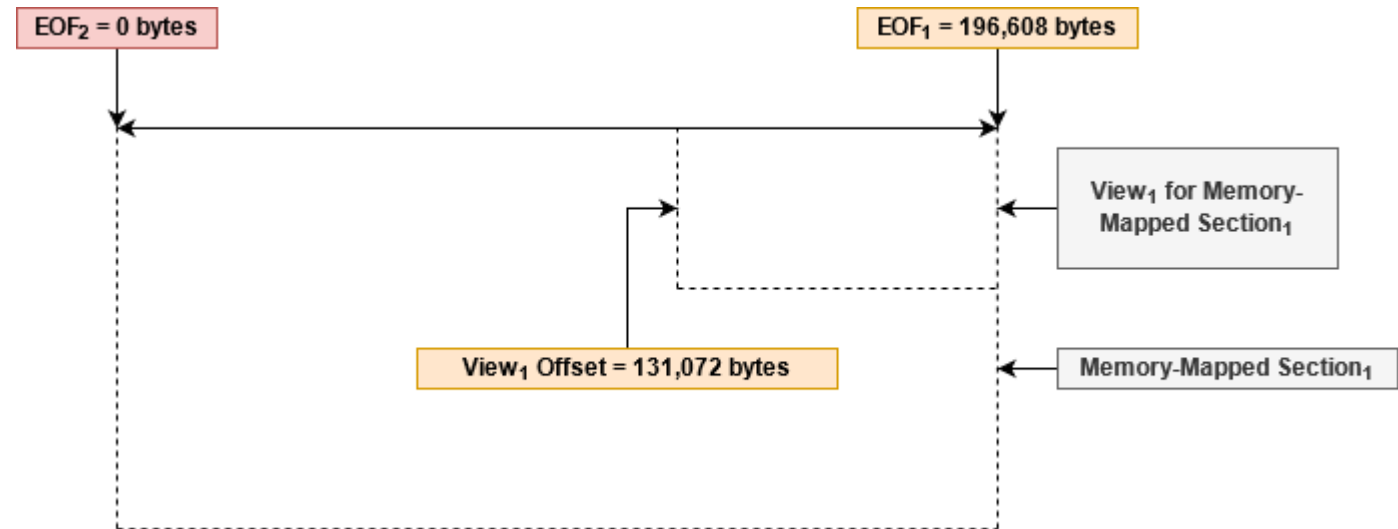
EOF₁ = 196,608 bytes

# Live-lock

- Overview of Issue
- Client$_1$ creates the following:
  - Memory-mapped section$_1$
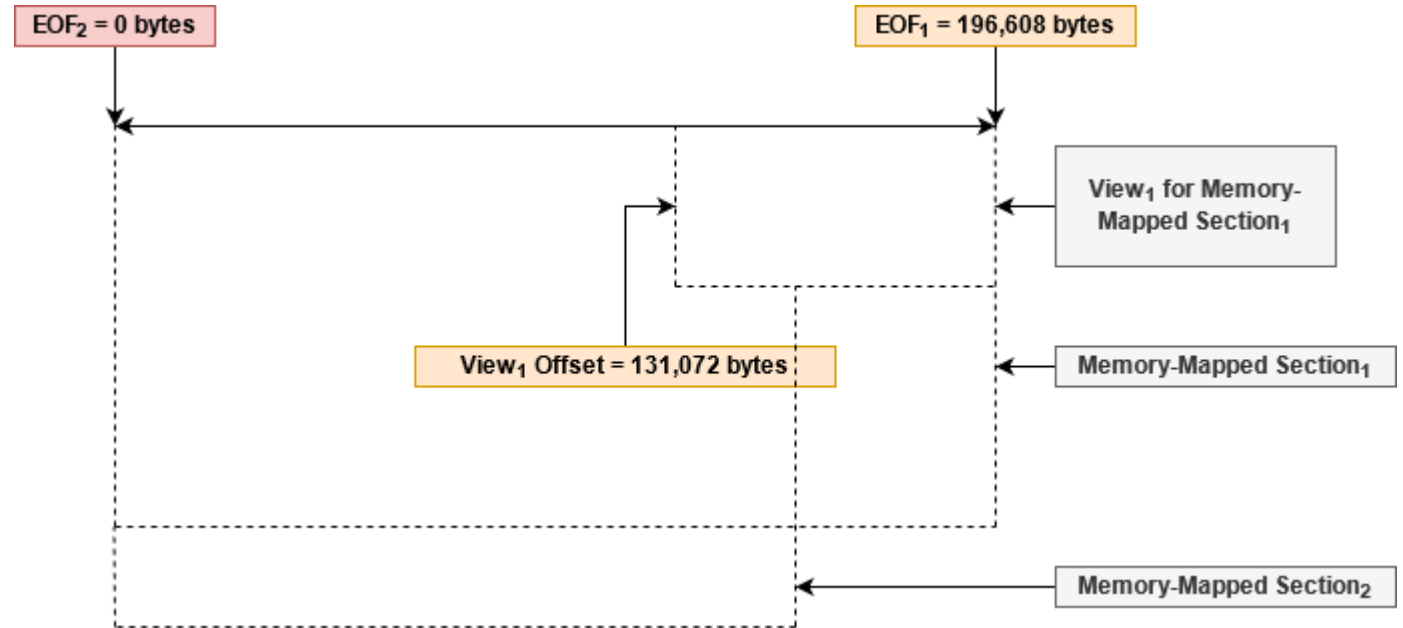  - View$_1$ with few dirty pages e.g., 16,384 bytes

# Live-lock

- Overview of Issue
- $Client_2$ does the following:
  - Open same file with shared read and write permissions.
  - Set the EOF = 0 (call this $EOF_2$)
- Oplock for $Client_1$ eventually downgrades from RWH → None
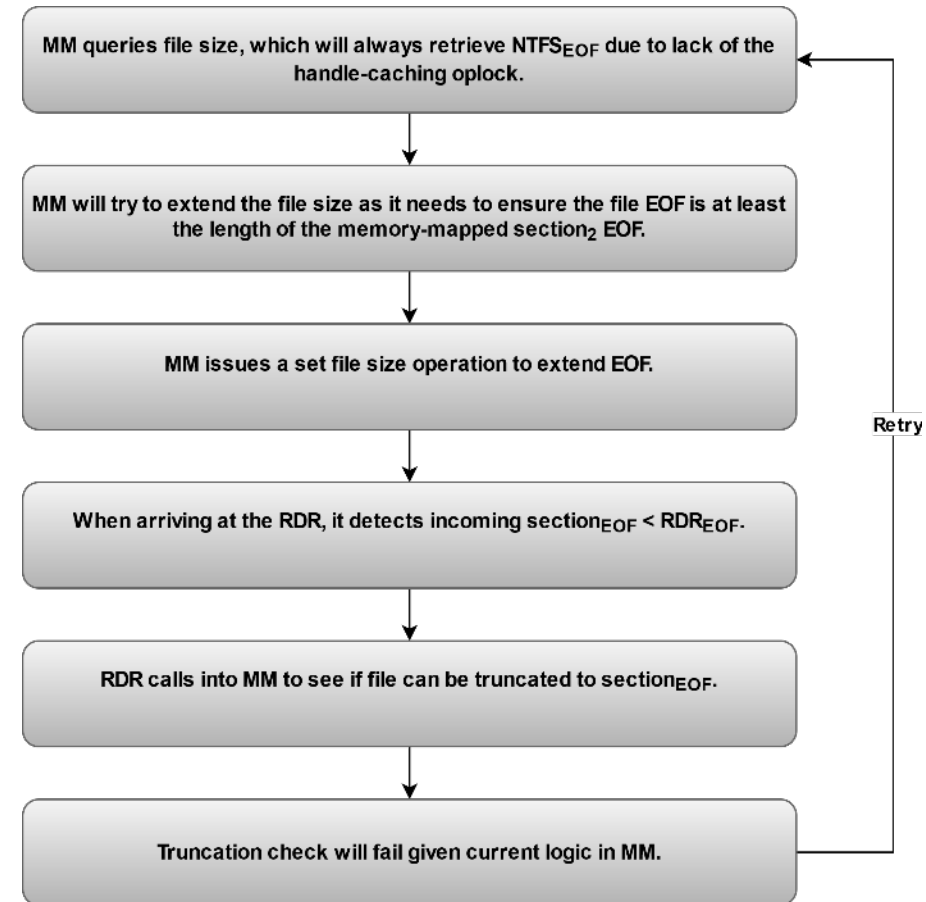- Dirty pages purged & flushed to down-level file system
- Data corruption occurs

# Live-lock

- Overview of Issue
- Client$_1$ creates memory-mapped section$_2$
  - **Triggers live-lock**

# Live-lock

- Overview of Issue
- <mark>$NTFS_{EOF}$ < Memory-Mapped Section #$2_{EOF}$ < $RDR_{EOF}$</mark>
  - Example earlier shows e.g.
    - 0 < 147,456 < 196,608
  - When the above condition is true, the endless loop occurs:

MM queries file size, which will always retrieve $NTFS_{EOF}$ due to lack of the handle-caching oplock.

↓

MM will try to extend the file size as it needs to ensure the file EOF is at least the length of the memory-mapped section$_2$ EOF.

↓

MM issues a set file size operation to extend EOF.

↓

When arriving at the RDR, it detects incoming section$_{EOF}$ < $RDR_{EOF}$.

↓

RDR calls into MM to see if file can be truncated to section$_{EOF}$.

↓

Truncation check will fail given current logic in MM.

Retry

# Live-lock

- Solution
  - Patching $RDR_{EOF}$ to the incoming section's EOF size
  - Preventing data loss in certain scenarios

# Please take a moment to rate this session.

Your feedback is important to us.

STORAGE DEVELOPER CONFERENCE

SDC 21