# Remote Access to Ultra-Low-Latency Storage

**Tom Talpey**
**Microsoft**

# Outline

- Problem Statement
- RDMA Storage Protocols Today
- Sources of Latency
- RDMA Storage Protocols Extended
- Other Protocols Needed

# Related SDC2015 Talks

- Monday – Neal Christiansen
- Tuesday – Jim Pinkerton, Andy Rudoff, Doug Voigt
- Wednesday – Chet Douglas
- Thursday – Paul von Behren

3

# Problem Statement

# RDMA-Aware Storage Protocols

☐ Focus of this talk – Enterprise / Private Cloud-capable storage protocols

  ☐ Scalable, manageable, broadly deployed

☐ SMB3 with SMB Direct

☐ NFS/RDMA

☐ iSER

☐ Many others exist

  ☐ Including NVM Fabrics, but not the focus here

5

# New Storage Technologies Emerging

- Advanced block devices
  - I/O bus-attached: PCIe, SSD, NVMe, …
  - Block or future Byte addressable
- Storage Class Memory ("PM" Persistent Memory)
  - Memory bus attached NVDIMM, …
    - Block or Byte accessible
  - Emerging persistent memory technologies
    - 3D XPoint, PCM, …
    - In various form factors

# Storage Latencies Decreasing

- Write latencies of storage protocols (e.g. SMB3) today down to 30-50us on RDMA
  - Good match to HDD/SSD
  - Stretch match to NVMe
  - PM, not so much ☺
- Storage workloads are traditionally highly parallel
  - Latencies are mitigated
- But workloads are changing:
  - Write replication adds a latency hop
  - Write latency critical to reduce

| Technology | Latency (high) | Latency (low) | IOPS |
|---|---|---|---|
| HDD | 10 msec | 1 msec | 100 |
| SSD | 1 msec | 100 μsec | 100K |
| NVMe | 100 μsec | 10 μsec (or better) | 500K+ |
| PM | < 1 μsec | (~ memory speed) | BW/size (>>1M/DIMM) |

Orders of magnitude decrease

7

**SDC** 15

# New Latency-Sensitive Workloads

❏ Writes!

   ❏ Small, random

      ❏ Virtualization, Enterprise applications

   ❏ MUST be replicated and durable

      ❏ A single write creates multiple network writes
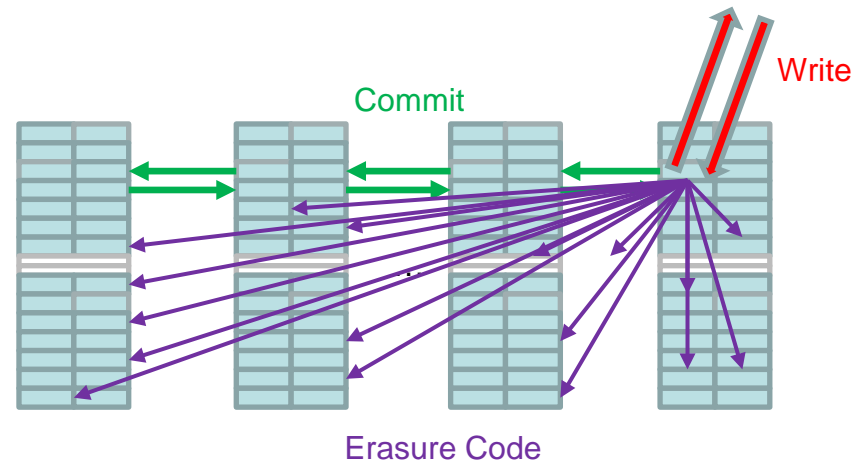
❏ Reads

   ❏ Small, random are latency sensitive

   ❏ Large, more forgiving

      ❏ But recovery/rebuild are interesting/important

8

# Writes, Replication, Network

- Writes (with possible erasure coding) greatly multiplies network I/O demand
  - The "2-hop" issue
- All such copies must be made durable before responding
  - Therefore, latency is critical!

Write

Commit

Erasure Code

# APIs and Latency

- APIs also shift the latency requirement
- Traditional Block and File are often parallel
- Memory Mapped and PM-Aware APIs not so much
  - Effectively a Load/Store expectation
  - Memory latency, with possibly expensive Commit
  - Local caches can improve Read (load) but not Write (store/remotely durable)
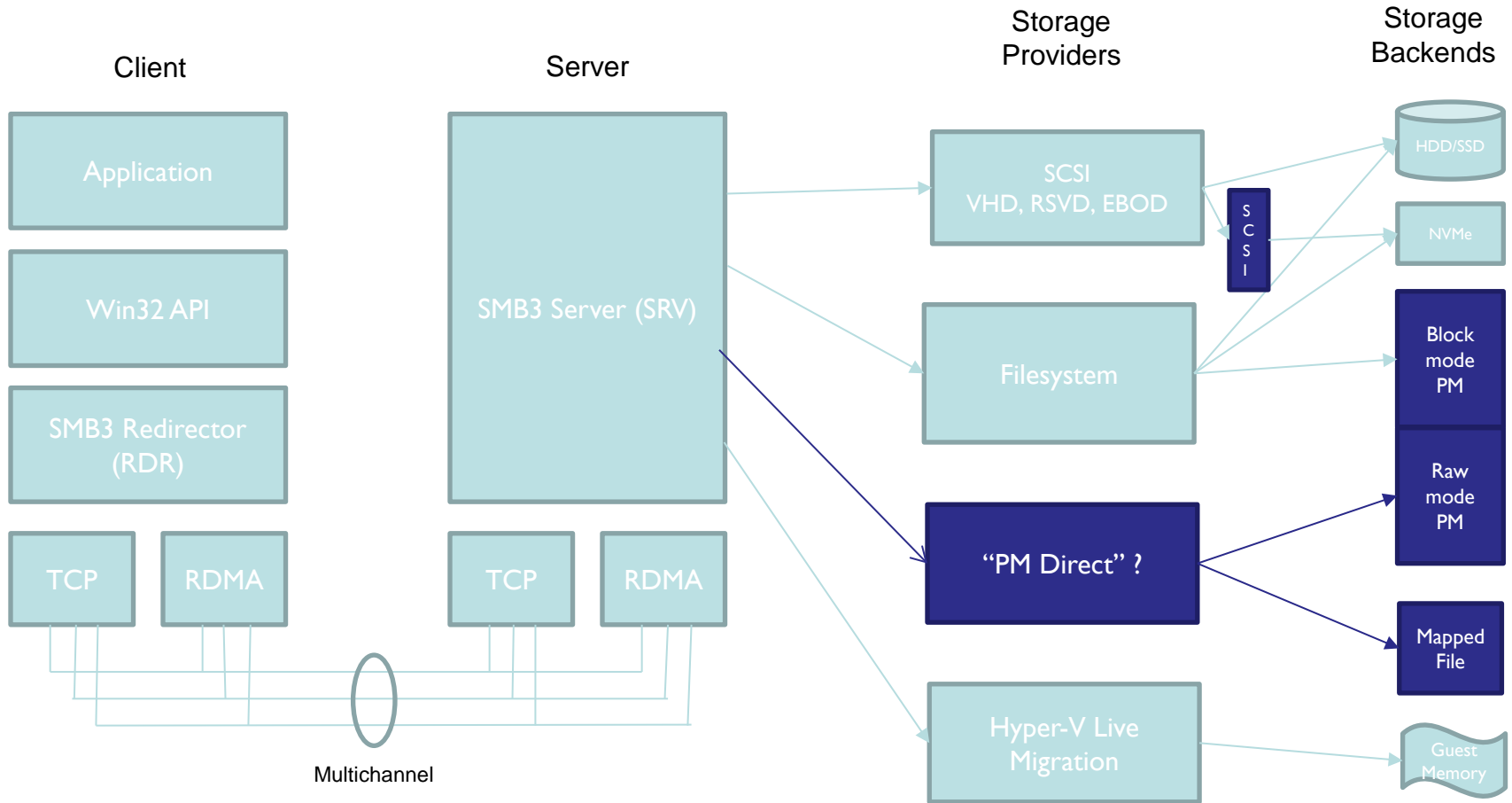
# RDMA Storage Protocols Today

# Many Layers Are Involved

- Storage layers
  - SMB3 and SMB Direct
  - NFS, pNFS and NFS/RDMA
  - iSCSI and iSER
- RDMA Layers
  - iWARP
  - RoCE, RoCEv2
  - InfiniBand
- I/O Busses
  - Storage (Filesystem, Block e.g. SCSI, SATA, SAS, …)
  - PCIe
  - Memory

# SMB3 Architecture (shameless plug)

- Principal Windows remote filesharing protocol
- Also an authenticated, secure, multichannel, RDMA-capable session layer
- Transport for
  - File system operations (REFS, NTFS, etc)
  - Block operations (VHDX, RSVD, "EBOD")
  - Hyper-V Live Migration (VM memory)
  - RPC (Named Pipes)
- Future transport for
  - Backend NVMe storage
  - Persistent Memory

# SMB3 Components (example)



Client

Application

Win32 API

SMB3 Redirector (RDR)

TCP    RDMA

Multichannel

Server

SMB3 Server (SRV)

TCP    RDMA

Storage Providers

SCSI
VHD, RSVD, EBOD

S C S I

Filesystem

"PM Direct" ?

Hyper-V Live Migration

Storage Backends

HDD/SSD

NVMe

Block mode PM

Raw mode PM

Mapped File

Guest Memory

SDC 15
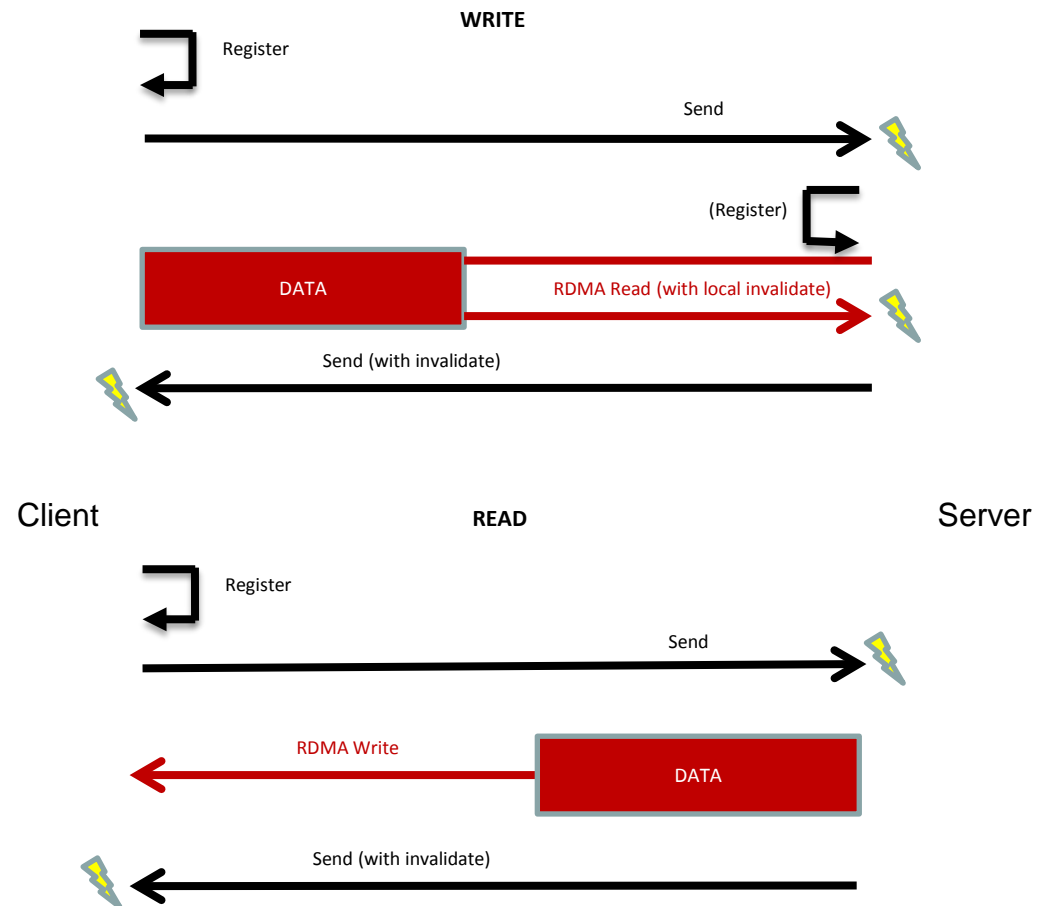
# Contributors to Latency

# RDMA Transfers – Storage Protocols Today

- ❑ **Direct placement model (simplified and optimized)**
  - ❑ Client advertises RDMA region in scatter/gather list
  - ❑ Server performs all RDMA
    - ❑ More secure: client does not access server's memory
    - ❑ More scalable: server does not preallocate to client
    - ❑ Faster: for parallel (typical) storage workloads
  - ❑ SMB3 uses for READ and WRITE
    - ❑ Server ensures durability
    - ❑ NFS/RDMA, iSER similar
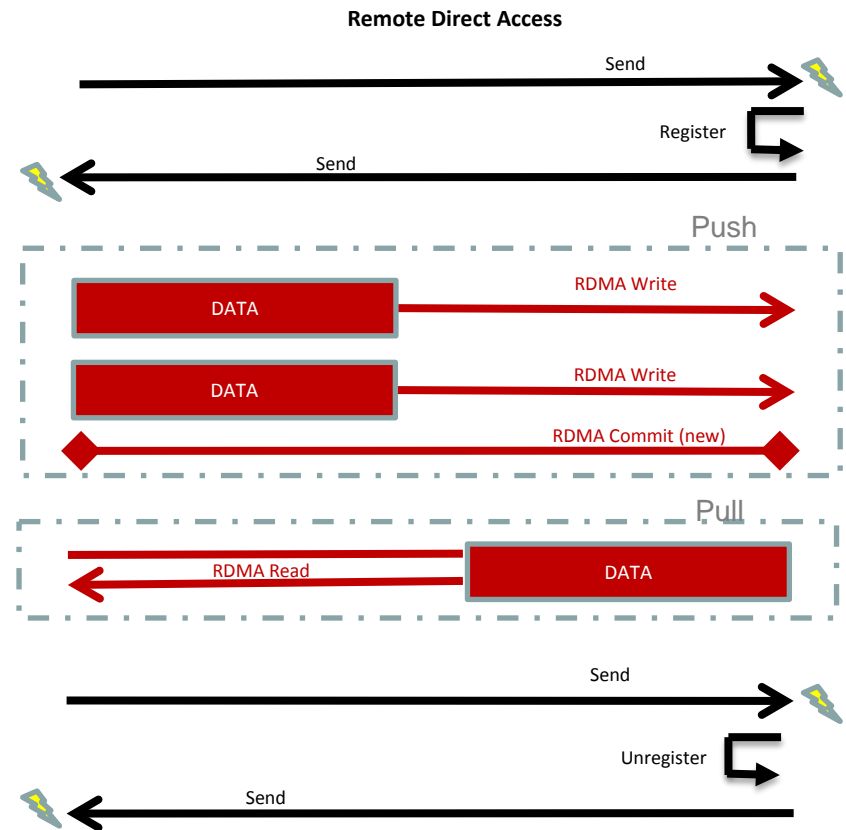  - ❑ Interrupts and CPU on both sides

**WRITE**

Register

Send

(Register)

DATA

RDMA Read (with local invalidate)

Send (with invalidate)

Client                    **READ**                    Server

Register

Send

RDMA Write

DATA

Send (with invalidate)

# Latencies

- Undesirable latency contributions
  - Interrupts, work requests
    - Server request processing
    - Server-side RDMA handling
  - CPU processing time
    - Request processing
  - I/O stack processing and buffer management
    - To "traditional" storage subsystems
  - Data copies
- Can we reduce or remove all of the above to PM?

# RDMA Storage Protocols Extended

# Push Mode (Schematic)

- Enhanced direct placement model
  - Client requests server resource of file, memory region, etc
    - MAP_REMOTE_REGION(offset, length, mode r/w)
  - Server pins/registers/advertises RDMA handle for region
  - <u>Client</u> performs all RDMA
    - RDMA Write to region
    - RDMA Read from region ("Pull mode")
    - No requests of server (no server CPU/interrupt)
      - Achieves near-wire latencies
  - **Client remotely commits to PM (new RDMA operation!)**
    - **Ideally, no server CPU interaction**
    - **RDMA NIC optionally signals server CPU**
    - **Operation completes at client only when remote durability is guaranteed**
- Client periodically updates server via master protocol
  - E.g. file change, timestamps, other metadata
- Server can call back to client
  - To recall, revoke, manage resources, etc
- Client signals server (closes) when done

**Remote Direct Access**

Send

Register

Send

Push

DATA — RDMA Write

DATA — RDMA Write

RDMA Commit (new)

Pull

RDMA Read — DATA

Send

Unregister

Send

19

# Storage Protocol Extensions

- ☐ SMB3
- ☐ NFSv4.x
- ☐ iSER

# SMB3 Push Mode (hypothetical)

- Setup – a new create context or FSCTL
  - Server registers and advertises w/r file by Handle
    - Or, directly to a region of PM or NVMe-style device!
  - Takes a Lease or lease-like ownership
- Write, Read – RDMA access by client
  - Client writes and reads directly via RDMA
- Commit – Client requests durability
  - Perform Commit, via RDMA with optional server processing
  - SMB_FLUSH-like processing for signaling if needed/desired
- Callback – Server manages client access
  - Similar to current oplock/lease break
- Finish – Client access complete
  - SMB_CLOSE, or lease manipulation

21

# NFS/RDMA Push Mode (hypothetical)

- Setup – a new NFSv4.x Operation
  - Server registers and advertises w/r file or region by filehandle
  - Offers Delegation or…
  - Via pNFS layout? (!)
- Write, Read – RDMA access by client
  - Client writes and reads via RDMA
- Commit – Client requests durability
  - Perform Commit, via RDMA with optional server processing
  - NFS4_COMMIT-like processing for signaling if needed/desired
- Callback – via backchannel
  - Similar to current delegation or layout recall
- Finish
  - NFS4_CLOSE, or delegreturn or layoutreturn (if pNFS)

22

SDC 15

# iSCSI (iSER) Push Mode (very hypothetical)

- ❑ Setup – a new iSER operation
  - ❑ Target registers and advertises w/r buffer(s)
- ❑ Write – a new <u>Unsolicited</u> SCSI-In operation
  - ❑ Implement RDMA Write within initiator to target buffer
    - ❑ No Target R2T processing
- ❑ Read – a new <u>Unsolicited</u> SCSI-Out operation
  - ❑ Implement RDMA Read within initiator from target buffer
    - ❑ No Target R2T processing
- ❑ Commit – a new iSER / modified iSCSI operation
  - ❑ Perform Commit, via RDMA with optional Target processing
  - ❑ Leverage FUA processing for signaling if needed/desired
- ❑ Callback – a new SCSI Unit Attention
  - ❑ ???
- ❑ Finish – a new iSER operation

SDC 15

# Other Protocols Extended

# RDMA protocols

❑ Need a remote guarantee of <u>Durability</u>
❑ RDMA Write alone is not sufficient for this semantic
  ❑ Completion at sender does not mean data was placed
    ❑ NOT that it was even sent on the wire, much less received
    ❑ Some RNICs give stronger guarantees, but <u>never</u> that data was stored remotely
  ❑ Processing at receiver means only that data was accepted
    ❑ NOT that it was sent on the bus
    ❑ Segments can be reordered, by the wire or the bus
    ❑ Only an RDMA completion at receiver guarantees placement
      ❑ And placement != commit/durable
  ❑ No Commit operation
❑ Certain platform-specific guarantees can be made
  ❑ But client cannot know them
  ❑ See Chet's presentation later today!

# RDMA protocol extension

- ❑ Two "obvious" possibilities
  - ❑ RDMA Write with placement acknowledgement
    - ❑ Advantage: simple API – set a "push bit"
    - ❑ Disadvantage: significantly changes RDMA Write semantic, data path (flow control, buffering, completion)
    - ❑ Requires significant changes to RDMA Write hardware design
      - ❑ And also to initiator work request model (flow controlled RDMA Writes would block the send work queue)
    - ❑ Undesirable
  - ❑ RDMA "Commit"
    - ❑ New operation, flow controlled/acknowledged like RDMA Read or Atomic
    - ❑ Disadvantage: new operation
    - ❑ Advantage: simple API – "flush", operates on one or more STags/regions (allows batching), preserves existing RDMA Write semantic (minimizing RNIC implementation change)
    - ❑ Desirable

26

# RDMA Commit (concept)

- RDMA Commit
  - New wire operation
  - Implementable in iWARP and IB/RoCE
- Initiating RNIC provides region list, other commit parameters
  - Under control of local API at client/initiator
- Receiving RNIC queues operation to proceed in-order
  - Like RDMA Read or Atomic processing currently
  - Subject to flow control and ordering
- RNIC pushes pending writes to targeted regions
  - If not tracking regions, then flushes all writes
- RNIC performs PM commit
  - Possibly interrupting CPU in current architectures
  - Future (highly desirable to avoid latency) perform via PCIe
- RNIC responds when durability is assured

# Local RDMA API extensions (concept)

- New platform-specific attributes to RDMA registration
  - To allow them to be processed at the server *only*
  - No client knowledge – ensures future interop
- New local PM memory registration
  - Register(region[], **PMType**, **mode**)
    - **PMType** includes type of PM
      - i.e. plain RAM, "commit required", PCIe-resident, any other <u>local</u> platform-specific processing
    - **Mode** includes disposition of data
      - Read and/or write
      - Cacheable after operation
    - Resulting handle sent by peer Commit, to be processed in receiving NIC under control of original Mode

28

# PCI Protocol Extension

- PCI extension to support Commit
    - To Memory, CPU, PCI Root, PM device, PCIe device, …
    - Avoids CPU interaction
    - Supports strong data consistency model
- Performs equivalent of:
    - CLFLUSHOPT (region list)
    - PCOMMIT
    - (See Chet's presentation!)

# Expected Goal

# Latencies

- Single-digit microsecond remote Write+Commit
    - (See Chet's presentation for estimate details)
    - Push mode minimal write latencies (2-3us + data wire time)
    - Commit time NIC-managed and platform+payload dependent
- Remote Read also possible
    - Roughly same latency as write, but without commit
- No server interrupt
    - Once RDMA and PCIe extensions in place
- Single client interrupt
    - Moderation and batching can reduce when pipelining
- Deep parallelism with Multichannel and flow control management

# Open questions

- Getting to the right semantic?
  - Discussion in multiple standards groups (PCI, RDMA, Storage, …)
  - How to coordinate these discussions?
  - Implementation in hardware ecosystem
  - Drive consensus from upper layers down to lower layers!
- What about new API semantics?
  - Does NVML add new requirements?
  - What about PM-aware filesystems (DAX/DAS)?
- Other semantics – or are they Upper Layer issues?
  - Authentication?
  - Integrity/Encryption?
  - Virtualization?

# Discussion?

**SDC** 15