

### The role of a InfiniBand and automated data tiering in achieving extreme storage performance

Cynthia Mcquire Oracle

### Introduction to InfiniBand



- Standard interconnect
  - Defined by the InfiniBand Trade Association
- Defined to facilitate low cost and high performance implementation
- Channel based I/O
- □ I/O consolidation
  - Communication, computation, management and storage over a single fabric



Number of	Per Lane Bandwidth				
IB Lanes	SDR 2.5Gb/s	DDR 5Gb/s	QDR 10Gb/s		
4X	10Gb/s	20Gb/s	40Gb/s		
8X	20Gb/s	40Gb/s	80Gb/s		
12X	30Gb/s	60Gb/s	120Gb/s		

### InfiniBand Feature Highlights

- Switch fabric links
  - I0Gb/s to 40Gb/s HCA links
  - Up to 120Gb/s switch-switch
- □ HW transport protocol
  - Reliable and unreliable
  - Connected and datagram
- □ Kernel bypass
  - Memory translation and protection tables
- Memory exposed to remote node
  - RDMA-read and RDMA-write
- Quality Of service
  - IO channels at the adapter level
  - Virtual lane at the link level
- □ Scalability and flexibility
  - Up to 48K nodes in subnet, up to 2<sup>128</sup> in network

- Network partitioning
  - Multiple networks on a single wire
- □ Reliable, lossless, self-managing fabric
  - End to end flow control
  - Link level flow control
  - Multicast support
  - Congestion control
  - Automatic path migration
- I/O Virtualization with channel architecture
  - Dedicated services to guests OSes
  - HW assisted protection and inter-process isolation
  - Enable I/O consolidation

### **QDR Product Portfolio**





Database – Oracle RAC

Low latency messaging

**ZFS** Storage Appliances

#### InfiniBand Network

- Redundant 40Gb/s switches
- Unified server & storage network





**Exadata** 

**Extreme Database Performance** 

# **Applications**

**Exadata** 

Exalogic





### **Solaris InfiniBand Overview**





SDC STORAGE DEVELOPER CONFERENCE

Application	IPE	Based Sockets Vario	us Block Clustered	Access to	SA	Subnet Administrator	
Levei	Diag Tools SM	cess Access MPI	s Access (Oracle RAC	C) Files	MAD	Management Datagram	
	User Level MAD	]			SMA	Subnet Manager Agent	
User APIs	OpenFabrics U	ser Level Verbs		uDAPL	PMA	Performance Manager Agent	
		SDP Lib		User Space	IPolB	IP over InfiniBand	
Upper Layer				Kernel Space	SDP	Sockets Direct Protocol	
Protocol	IPoIB CM UD	EolB SDP SRP i	SER FCoIB RDS rNFS	Lustre LND	SRP	SCSI RDMA Protocol (Target & Initiator)	
		• • • • • • • • • • • • • • • • • • •			iSER	iSCSI RDMA Protocol (Target & Initiator)	
Mid-Layer	ass se			ss	RDS	Reliable Datagram Service	
	AG Client MAD SMA	UDAPL	User Direct Access Programming Lib				
	InfiniBand Transport Framework					RDMA for RPC under NFS	
					OpenSM	Subnet Manager	
Drivers	Hermo	n	Τανο	Key	DONE Apps & Access		
Hardware	ConnectX	ConnectX-2	Arbei Mem-Full Tavor		L	InProgress Methods for using MNX HCA Stack	



- Primary vendor of HCA silicon is Mellanox
- **Two drivers supported:** 
  - tavor(7D) for InfiniHost (PCI-X), InfiniHost III Ex (PCIe), up to DDR rates
  - hermon(7D) for ConnectX[-2] (PCIe gen 1 & 2), up to QDR rates
- Hermon (7D) offers additional functional and performance features
  - 256 MSIx vectors used to balance CPU loading
  - Offloads LSO (network), FRWR (memory registration)
  - Capabilities Relaxed ordering support, HCA hot-plug
- □ Firmware update via fwflash(1m)
- □ Integration with Solaris FMA
- □ Show picture of HCA from Mellanox PRM, qps. cqs

### InfiniBand Stack: Transport (IBTF)



- Enables multiple services over IB transport
  - Compliant with 1.2 specification
  - Utilizes network services to map IP address to IB address
- **Exports** Transport Interface (TI) upward to service drivers
  - The "Verbs" interface
- **Requests services from HCA Driver** 
  - Channel Interface (CI) downward to HCAs
- **Provides Management Datagram (MAD) interface to components** 
  - Communication manager establishes connections through fabric
  - Device manager puts attached I/O devices in Solaris devfs
  - Agents (SMA, PMA, etc), Apps (OpenSM) can access fabric components in standard ways

### InfiniBand Stack: IP over IB



### ETF standard compliant

- 4391 IP transmission over IB
- 4755 IP transmission over the RC and UC modes of IBA to support large MTUs
- 4390 DHCP over IB
- 20 byte link layer address
  - Most common "porting" issue
  - Does not persist across reboot
- □ Maps IP subnet to IB partition
- Enables IP addressing for other IB ULPs



# InfiniBand Stack: IPoIB (continued)

- TCP checksum offload and LSO, interrupt moderation supported
- □ dladm(1m) administration
  - "link-mode" property to specify
    "cm" or "ud" data transfer mode
  - IB "vnics" through new \*-part subcommands
- □ IPMP for high availability
- Interoperable with Linux and Windows
- Observability
  - Snoop and Wireshark for IB



SD

SNIA SANTA CLARA, 2011

#### **Ethernet over IB (EoIB)**

- Protocol developed by Mellanox and used with the Sun NanoMagnum Gateway (NM2-GW) for bridging IB and Ethernet networks
  - Protocol defines the procedure for tunneling both unicast and multicast Ethernet packets over IB using the Unreliable Datagram (UD) transport
    - The Ethernet packet includes the standard header, VLAN tags, if appropriate, and the payload.
  - Protocol defines control messages used by Gateway management software to associate IB connected servers with Ethernet ports on Gateway
- EoIB protocol server-side endpoint on Solaris is implemented as standard network driver (NIC)
  - □ EoIB NICs behave like regular Ethernet NICs
    - Networking stack works seemlessly over EoIB NICs

### **IB Connected Storage: iSER**



#### □ iSCSI extensions to RDMA

- Defined by the IETF as RFC 5046 with extensions defined by IBTA for InfiniBand
- iSER extends the data transfer model of iSCSI
  - □ Reliable protocol and data integrity by IB
  - □ Zero copy using RDMA
  - □ Minimal CPU overhead per IO because transport is IB
- iSER still maintains compliance with iSCSI
  - □ Uses exactly the same iSCSI PDUs
  - Utilizes existing iSCSI infrastructure (e.g. bootstrapping, MIB, negotiation, naming and discovery, and security)
  - Utilizes an iSCSI mechanism, "login key negotiation" to determine whether to use the iSCSI or iSER data transfer models
  - □ Add a 12 Byte iSER header before iSCSI PDU (for control and buffer advertisement)

### **IB Connected Storage: iSER**





#### **Example SCSI Read**

- □ Initiator Send Command PDU (Protocol data unit) to Target
- □ Target return data using RDMA Write
- Target send Response PDU back when completed transaction
- □ Initiator receives Response and complete SCSI operation

### **IB Connected Storage**



#### □ iSER: iSCSI Extensions for RDMA

- itadm(1m) and iscsiadm(1m) used to setup initiator
- Target support in Solaris 11 Express
- Except for data movement, iSER interoperates as an iSCSI Initiator and Target
  - If both the iSER initiator and target are unable to negotiate an RDMA channel, or the channel fails and it can not be re-established, data movement regress back to TCP/IP

#### SRP: SCSI RDMA Protocol

- ANSI T-10 standard
- Uses RDMA for high BW, low latency, low overhead
- First IB storage protocol implemented
- Considered the IB industry standard
  - □ Initiator supported on VMware, Linux and Windows
- Target support in Solaris 11 Express
- Uses IB addressing and device management model

# **IB Connected Storage: NFS and IB**



- □ NFS and InfiniBand (NFS/RDMA)
- □ Semantically, NFS is NFS
  - Applications see no difference in semantics
  - RDMA helps in throughput, latency, efficiency
- RDMA under NFSv3 and NFSv4 protocols
- Use Reliable Connected
- □ NFS runs on standard TCP and UDP networks
  - NFS-over-IPoIB requires no special setup
  - NFS-over-IPoIB does benefit from faster media
- Remote DMA (RDMA) on Infiniband
  - RDMA is transparently negotiated if available
- □ mount(IM) used to set "proto=tcp" or "proto=rdma"
  - Default mode will negotiate
- **RDMA** interop tested continuously at Connectathon / Bakeathon events
- **Tested to interoperate with Linux** 
  - Requires at least 2.6.31 kernel

### **NFS/RDMA (continued)**



#### **NFS** reads

- Client posts a write chunk list
- □ Server transfers the data to the client using RDMA\_WRITE
- □ Server notifies the client with inline reply



### **The Converged Fabric Stack**





### Not there yet...





#### Benefits of InfiniBand negated by bottlenecks of the back-end

- Filesystem
- IO bus
- HBS
- Storage fabric
- Disks

2011 Storage Developer Conference.  $\ensuremath{\mathbb{C}}$  Oracle Corporation. All Rights Reserved.

### **Storage Stack for High Performance**





- □ Kernel level Common SCSI Target COMSTAR
- **Reduced copy between filesystem and protocol layers** 
  - Shared buffers with ZFS and NFS/iSCSI/iSER
- **ZFS -** Pooled storage: converged filesystems and volumes
  - Pipelined I/O
  - Dynamic striping across
    - Distributes load across devices
  - Intelligent prefetch
    - Multiple independent prefetch streams
  - Variable block sizes
    - □ Large blocks: less metadata, higher bandwidth
    - □ Small blocks: more space-efficient for small objects
    - Record-structured file (e.g. databases) sizes matched to filesystem match to avoid read/modify/write

### **Storage Hardware for High Performance**

SDC STORAGE DEVELOPER CONFERENCE



### Hybrid Storage Pool – ZFS





- □ Disks for high capacity, high bandwidth, low power
  - Great aggregate throughput
- **DRAM** as primary ZFS cache
  - Used as very fast adaptive replacement cache (ARC)
- Logzilla: write-optimized flash
  - Used as the ZFS intent log (ZIL)
  - Synchronous writes don't wait for disk; acknowledged as soon as they reach the logzilla
  - Data asynchornously streamed to disk
- Readzilla: read optimized flash
  - Used as second tier cache (L2ARC)



### **Solaris**

- Integrated into existing command set: ifconfig, dladm, itadm, iscsiadm
- OFVA verbs library to 'roll your' diagnostic tools
- Wireshark, snoop

### **Linux**

- Wireshark
- OFED stack
- Commands, libraries, diagnostic, performance tool

### **Administrative Ease**



Sun CRACLE	SUN ZFS STORAGE	5 7320				Super-Use	en@auto7320-08_Loc	OUT HELP
								nalytics
		SER	VICES STORAGE	NETWORK	SAN CI	USTER USERS	PREFERENCES	ALERTS
Network		Network Data	link		CANCE	APPLY	Addresses	Routing
To configure net an object to view	working, build Datalinks its relationship to othe	Name ipoib-1					REVERT	APPLY
Devices	6 total	<b>.</b>						4 total
BUILT-IN		Status					yin inb0	L Ū
🗰 igb0	link down	Proportion			IR Partition		, via iguo	1 1
igb1 🞆	1Gb (full)	Fropences				1	pffff_ibp0	<i>и</i> . ш
酬 igb2	link down		Partition	Key		_	uis intri	L T
tigb3	link down		Link M	ode Connect	ed Mode	\$	s, vialigui	e tit
PCIe 0							pffff_ibp1	ar. w
<b>≣</b> €ibp0	32Gb (port 1)	Partition Dev	rices 2/2 available	LACP A	Aggregation			
🛲 ibp1	32Gb (port 2)		0x212800013e7c4t			32Gb (port 1)		
		💛 🛲 ibp1	0x212800013e7c50			32Gb (port 2)		

- Configured just like IP datalinks and interfaces over an IB port
- Specify the partition key and done!

 Configured just like iSCSI targets over IB network interface

SUN ZFS STORAGE	7320		3888	22	201	Super-User	@auto7320-08 LO	GOUT HELP
	SERVICES	STORAGE	NETWORK	SAN	CLUSTER	USERS	PREFERENCES	ALERTS
Storage Area Network (	New iSCSI Target			CA	NCEL	ок	Targets	Initiators
To share LUNs only via particular targe respectively. To create a group or add t		<b>T</b>	ON <b>O</b> • •				REVERT	
Fibre Channel Ports		Target	QN O Auto-as	ssign				
		A	lias iser-ipoib					
	Initiator authentication mode O None							
		-	O RADIU	S				
	Target CHAP name Target CHAP secret igb 1 igb 1							

### **Administrative Ease**





- Observability and diagnosability
  - CPU usage
  - Protocol IOPS
  - IPoIB and RDMA
    - □ Throughput
    - □ Latency



### search.oracle.com

ZFS Storage Appliances

