# SQLintersection

Monday, 9:00-4:30

## Analyzing and Improving I/O Subsystem Performance

Glenn Berry

glenn@sqlskills.com

SQL *intersection*

# Glenn Berry

- **Consultant/Trainer/Speaker/Author**
- **Principal Consultant, SQLskills.com**
  - Blog: http://www.SQLskills.com/blogs/Glenn
  - Twitter: @GlennAlanBerry
  - Regular presenter at worldwide conferences on hardware, scalability, and DMV queries
  - Author of SQL Server Hardware
  - Chapter author of Professional SQL Server 2012 Internals and Troubleshooting
  - Chapter author of MVP Deep Dives Volumes 1 and 2
- **Instructor-led training: Immersion Events**
- **Online training:** pluralsight    **http://pluralsight.com/**
- **Consulting: health checks, hardware, performance, upgrades**

SQL *intersection*

# Overview

- **Three main metrics for storage performance**
- **SQL Server I/O workload metrics**
- **Tools for testing storage subsystems**
- **Primary storage types for SQL Server**
- **RAID levels and SQL Server workloads**
- **Some comparative storage metrics**
- **Improving I/O performance**
- **Index and workload tuning**

# Three Main Metrics for Storage Performance

- **Latency (ms)**
- **Input/output operations per second (IOPS)**
- **Sequential throughput (MB/sec or GB/sec)**
  - These three measurements are all related, so you can't just look at one of them in isolation without knowing the others
  - Storage vendors tend to show their best-case numbers in isolation

# Latency

- **Latency is the time it takes for an I/O to complete**
  - Sometimes called response time or service time
- **Measurement starts when the OS sends a request to the drive (or controller) and ends when the drive finishes processing the request**
  - Reads are complete when the OS receives the data
  - Writes are complete when the drive informs the OS it has received the data
    - The data may still be in a DRAM cache on the drive or controller

**SQL**
*intersection*

# Input/Output Operations per Second

- **Input/output operations per second (IOPS)**
  - This metric is directly related to latency
    - Constant latency of 1ms means a drive can process 1,000 I/Os per second with a queue depth of 1
  - As more I/Os are added to the queue, latency will increase
  - Flash storage can read/write to multiple NAND channels in parallel
- **IOPS = Queue Depth/Latency**
- **IOPS by itself does not consider the transfer size**
  - You need to know the transfer size when looking at an IOPS measurement
  - You can translate IOPS to MB/s and MB/s to latency as long as you know the queue depth and transfer size

**SQL** *intersection*

# Sequential Throughput

- **Sequential throughput (MB/sec or GB/sec)**

- **MB/sec = IOPS * Transfer Size**
  - 556 MB/sec = 135,759 IOPS * 4096 bytes transfer size
  - 1112 MB/sec = 135,759 IOPS * 8192 bytes transfer size

- **Sequential throughput often gets short-changed in enterprise storage**
  - Bandwidth limitations from the storage interface directly affect this
  - 1Gbps iSCSI limited to about 100 MB/sec
  - 4Gbps FC limited to about 400 MB/sec

- **The disks may be so busy that they can't deliver full rated throughput**
  - This is fairly common with magnetic disks in DAS and shared SANs

SQL
*intersection*

# The Importance of Sequential Throughput

- **Sequential throughput is critical for many database server activities**
  - Full database backups and restores
    - Initializing AlwaysOn AG replicas, database mirrors, replication subscribers, log shipping secondaries
  - Index creation and rebuilds
    - Good sequential throughput makes it much easier to do index tuning
  - DW or reporting workload large sequential scans
    - Very important when data does not fit in the buffer pool
    - Unfortunately, SQL Server 2014 BPE does not help with sequential scans

**SQL**
*intersection*

# Demo

**Measuring sequential throughput**

SQL *intersection*

# **Important SQL Server I/O Workload Metrics (1)**

- **What is the read vs. write ratio of the workload?**
    - You can use my DMV Diagnostic Queries to determine this at the file level
    - Ratios will be different for different SQL Server file types and workloads
- **What are the typical I/O rates (IOPS and throughput)?**
    - Reads/sec, writes/sec (PerfMon) is IOPS
    - Disk read bytes/sec, disk write bytes/sec (PerfMon) is throughput

**SQL**
*intersection*

# Important SQL Server I/O Workload Metrics (2)

- **What is the average logical disk-level latency?**
  - Average disk sec/read, average disk sec/write (PerfMon) is latency
  - You can use one of my DMV Diagnostic Queries to determine this at the disk level
- **What is the average file-level I/O latency?**
  - You can see this in Windows Resource Monitor
  - You can use one of my DMV Diagnostic Queries to determine this at the file level

SQL
*intersection*

# Methods for Measuring I/O Performance (1)

- **Task Manager in Windows Server 2012 and newer**
  - Depending on what kind of storage you are using, not always useful
- **Disk section in Windows Resource Monitor**
  - Shows file-level response time in ms
- **Logical Disk Counters in Performance Monitor**
  - Shows disk-level metrics

SQL
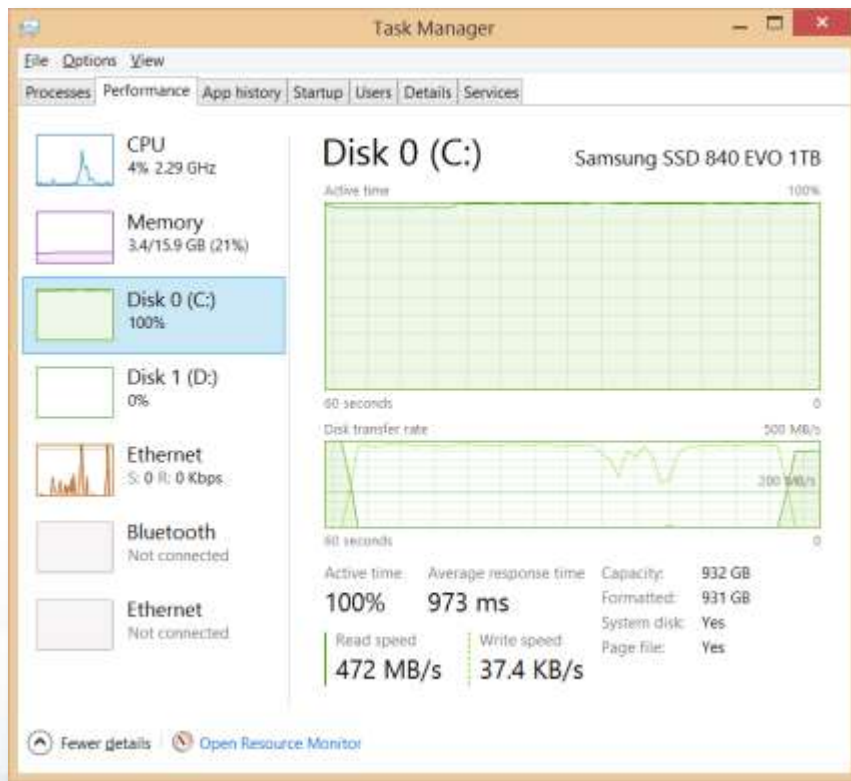*intersection*

# Methods for Measuring I/O Performance (2)

- **Disk Benchmark Tools**
  - CrystalDiskMark
    - http://bit.ly/1vm5dPe
  - Microsoft DiskSpd
    - http://bit.ly/1whNzQL
  - Microsoft SQLIO
    - SQLIO is deprecated, so I would not use it for new testing
    - http://bit.ly/1obVdIV
- **SQL Server DMV Diagnostic Queries**
  - http://bit.ly/Q5GAJU
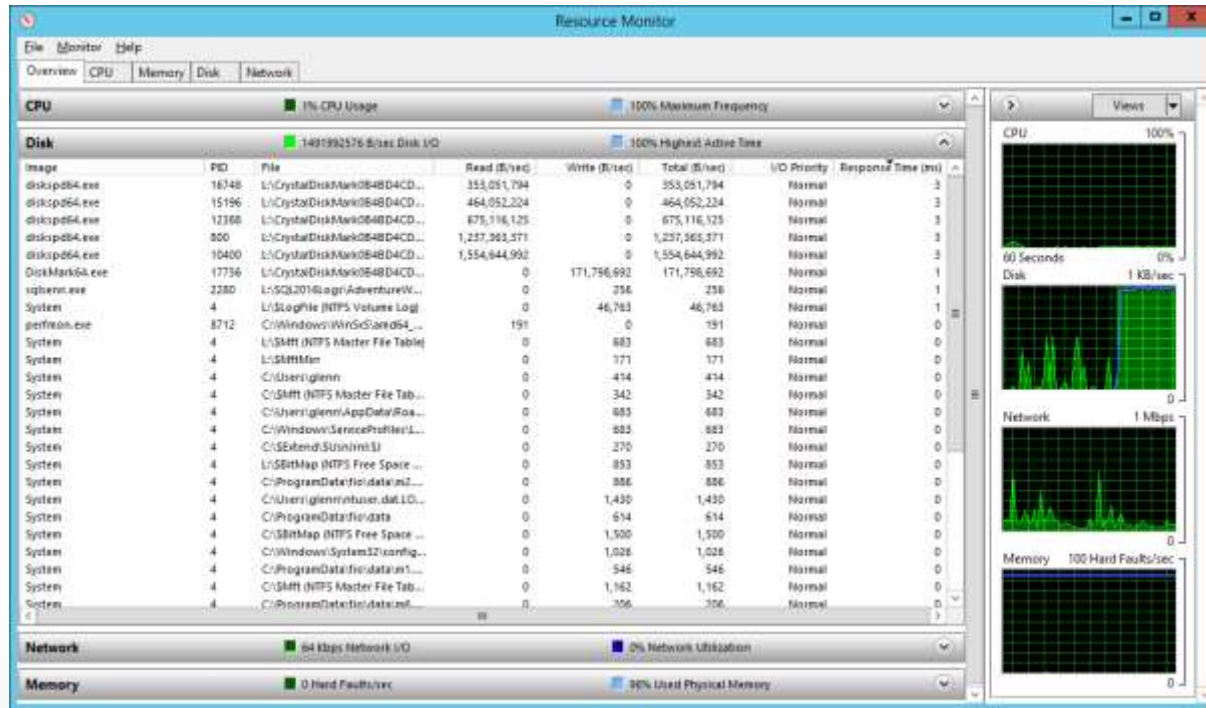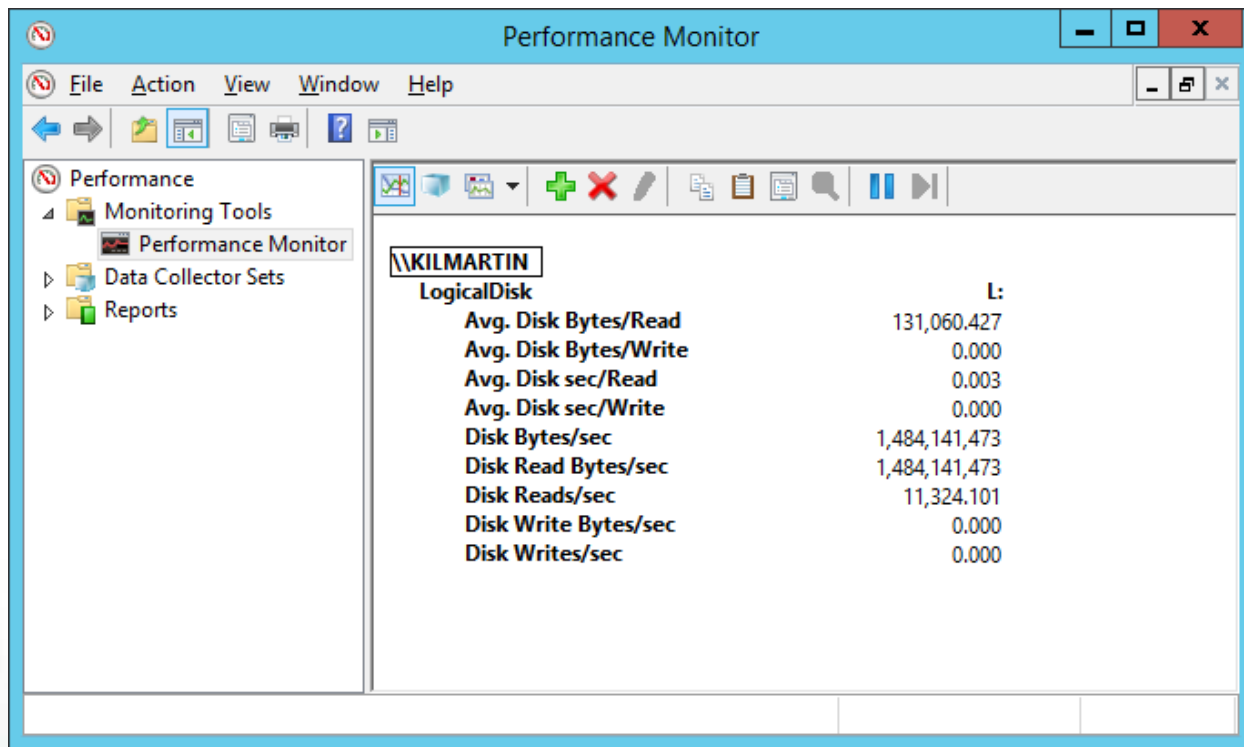
SQL
*intersection*

# Disk Performance in Windows Task Manager

# Disk Performance in Windows Resource Monitor

# LogicalDisk Counters in Performance Monitor

# Demo

**I/O related DMV queries**

(from SQL Server DMV Diagnostic Queries)

**SQL** *intersection*

# Common DMV I/O Query Result Patterns

- **Very common to see high write latency to tempdb data files**
  - Make sure you have multiple data files (nstart with 4-8) that are all the same size (follow Bob Ward's guidance)
  - Make sure you are using TF 1118 (not necessary with SQL Server 2016)
  - Consider using local flash-based storage for tempdb, if your workload needs it
- **Common to see high read latency from user database data files**
  - Look for signs of memory pressure, consider adding more RAM and doing standard workload and index tuning
  - Consider using SQL Server 2014 BPE (esp for Standard Edition) with OLTP workloads
    - Make sure to use fast, local flash storage for BPE file

**SQL**
*intersection*

# SQL Server 2014 Buffer Pool Extension (BPE)

- **Allows use of a cache file in the file system to cache clean pages**
  - Makes buffer pool appear to be larger
  - Can help with OLTP workloads with lots of random reads
  - Does not help very much with large sequential reads (by design)
  - Cache file must be at least as large as max server memory setting
  - Make sure to use fast, local flash storage!
- **Most useful for SQL Server 2014 Standard Edition**
  - SQL Server 2014 Standard Edition has 128GB RAM limit
  - BPE file can be 4X max server memory setting in Standard Edition
  - Can also be useful with small VMs

SQL
*intersection*
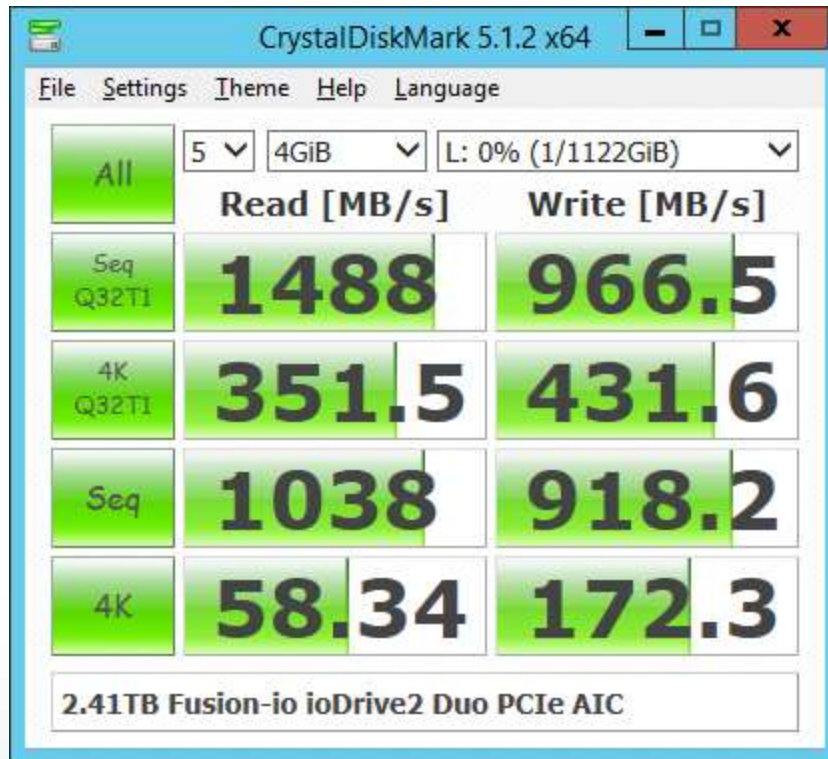
# Demo

**Using Buffer Pool Extensions**
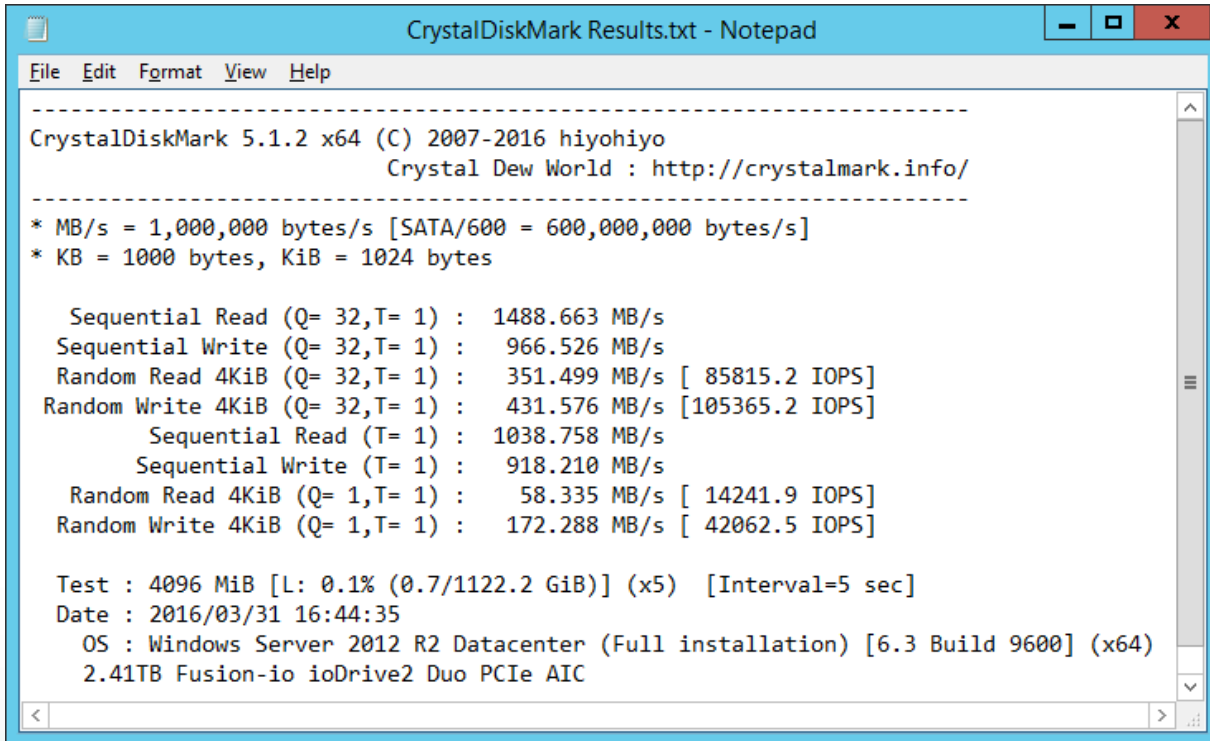
**SQL** *intersection*

# Using CrystalDiskMark To Test Your LUNs

- **Always use CrystalDiskMark for initial quick testing of logical drives**
  - CrystalDiskMark does not work with mount points
  - Test each logical drive before you install SQL Server if possible
  - Make sure to test with a large test file size
    - This will minimize the influence of any hardware cache
  - Make sure to test with both random and non-random test file types
    - Random data is not compressible, non-random data is compressible
    - Some controllers use write compression
  - Make sure to select at least five test runs for the test
    - This reduces the chances of outliers skewing the results

SQL
*intersection*

# CrystalDiskMark Results (Graphical)

# CrystalDiskMark Results (Text)



```
-------------------------------------------------------------------
CrystalDiskMark 5.1.2 x64 (C) 2007-2016 hiyohiyo
                      Crystal Dew World : http://crystalmark.info/
-------------------------------------------------------------------
* MB/s = 1,000,000 bytes/s [SATA/600 = 600,000,000 bytes/s]
* KB = 1000 bytes, KiB = 1024 bytes

   Sequential Read (Q= 32,T= 1) :  1488.663 MB/s
  Sequential Write (Q= 32,T= 1) :   966.526 MB/s
  Random Read 4KiB (Q= 32,T= 1) :   351.499 MB/s [ 85815.2 IOPS]
 Random Write 4KiB (Q= 32,T= 1) :   431.576 MB/s [105365.2 IOPS]
          Sequential Read (T= 1) :  1038.758 MB/s
         Sequential Write (T= 1) :   918.210 MB/s
   Random Read 4KiB (Q= 1,T= 1) :    58.335 MB/s [ 14241.9 IOPS]
  Random Write 4KiB (Q= 1,T= 1) :   172.288 MB/s [ 42062.5 IOPS]

  Test : 4096 MiB [L: 0.1% (0.7/1122.2 GiB)] (x5)  [Interval=5 sec]
  Date : 2016/03/31 16:44:35
    OS : Windows Server 2012 R2 Datacenter (Full installation) [6.3 Build 9600] (x64)
    2.41TB Fusion-io ioDrive2 Duo PCIe AIC
```

# Demo

**Using CrystalDiskMark 5.1.2**

SQL *intersection*

# Using SQLIO To Test Your Storage

- **SQLIO does not require or use SQL Server for its testing**
  - It simply allows you to stress your I/O subsystem in a fairly controlled manner
  - SQLIO is deprecated (but still works if you want to use it). I would not...
- **SQLIO has many configuration options**
  - Can be time consuming to run full suite of tests. Can be dangerous on shared SANs
- **You can use old style command prompt or PowerShell to run tests**
- **Reference:**
  - SQLIO, PowerShell and storage performance: measuring IOPs, throughput and latency for both local disks and SMB file shares
  - http://bit.ly/1n7jm0M

SQL *intersection*

# Using DiskSpd To Test Your Storage

- **DiskSpd does not require or use SQL Server for its testing**
  - It is a new tool from Microsoft and is far more flexible and powerful than SQLIO
  - It gathers much more information than SQLIO does
- **You can use old style command prompt or PowerShell to run tests**
- **Example command line:**
  - C:\DiskSpd\diskspd.exe -c100G -d10 -r -w0 -t8 -o8 -b8K -h -L X:\testfile.dat
- **Reference:**
  - DiskSpd, PowerShell and storage performance: measuring IOPs, throughput and latency for both local disks and SMB file shares
  - http://bit.ly/1CeQauw

**SQL**
*intersection*

# Demo

## Using Microsoft DiskSpd

**SQL** *intersection*

# Primary Storage Types for SQL Server

- **Several different storage types are commonly used**
  - Internal drives (3.5", 2.5", or 1.8")
  - Direct-attached storage (DAS)
  - Storage area networks (SAN)
  - PCIe flash-based AIC storage cards
  - Server Message Block (SMB) 3.0/3.02 file shares
  - Scale-Out File Servers (SOFS)
  - Storage Spaces Direct (S2D) in Windows Server 2016

**SQL**
*intersection*

# Internal Drives

- **Internal drives can be adequate for many workloads**
  - Possible to have up to (28) 2.5" drives in some two-socket servers
  - Lots of capacity and drive performance possible with this many drives
  - Can be very well suited for AlwaysOn AG node storage
- **Rack-mount server vertical form factor affects number of drive bays**
  - Drive size (3.5", 2.5", or 1.8") affects drive density
- **Use the best hardware RAID controller(s) available for your server model**
  - Premium RAID controllers have faster processors and larger cache sizes
  - This is especially important for parity-based RAID levels or flash storage

**SQL** *intersection*

# Direct-Attached Storage (DAS)

- **External storage enclosure with multiple drive bays**
  - Typically (14-24) 2.5" drives in a single external storage enclosure
  - Try to dedicate at least one RAID controller to each storage enclosure
  - Storage enclosures should have dual power supplies
- **DAS is easy to configure and manage**
  - Does not require special training or expertise or a cranky SAN administrator…
  - Does require planning and common sense
- **Can provide excellent sequential read/write performance**
  - Limited by PCIe slot bandwidth and RAID controller performance
  - Can be very well suited for AlwaysOn AG node storage

## SQL
*intersection*

# DAS Considerations

- **Use one dedicated RAID controller per storage enclosure**
  - You may even want two RAID controllers per enclosure
  - Use the best PCIe RAID controller available
  - Make sure the hardware cache is enabled
- **Pay attention to the PCIe slot throughput limits**
  - Does not require special training or expertise
  - Does require planning and common sense.
  - Try to dedicate the hardware RAID controller cache to writes
  - Disable read-ahead caching
  - The SQL Server buffer pool is a better read cache than the hardware RAID cache

**SQL**
*intersection*

# PCIe Slot Bandwidth Limits

- **PCIe 1.0 Bus (one-way)**
  - x4 slot: 750MB/sec
  - x8 slot: 1.5GB/sec
- **PCIe 2.0 Bus (one-way)**
  - x4 slot: 1.5-1.8GB/sec
  - x8 slot: 3.0-3.6GB/sec
- **PCIe 3.0 Bus (one-way)**
  - x4 slot: 3.0-3.6GB/sec
  - x8 slot: 6.0-7.2GB/sec
- **x4 and x8 refer to the number of lanes the slot supports**

**SQL** *intersection*

# PCIe Version Support

- **PCIe 3.0 Support**
  - Intel Xeon E5, E5 v2, E5 v3 and E5 v4 families
    - (Sandy Bridge-EP, Ivy Bridge-EP, Haswell-EP, and Broadwell-EP)
  - Intel Xeon E7 v2 and E7 v3 families (Ivy Bridge-EX and Haswell-EX)

- **PCIe 2.0 Support**
  - Intel Xeon E7 family (Nehalem-EX)
  - Older Intel processors
  - All current AMD processors

# PCIe Flash Storage

- **Flash-based storage on a PCIe expansion card**
  - Uses very high bandwidth PCIe slot instead of SAS/SATA port
  - New products using Non Volatile Memory Express (NVMe) have excellent performance!
  - Type and speed of PCIe slot can be limiting factor
  - PCIe flash storage cards can deliver extremely high I/O performance
    - Very high sequential throughput (up to 6.7GB/sec)
    - Extremely high random I/O performance (up to 1.3 million IOPS)
  - Flash storage cards use less electrical power than multiple magnetic drives
    - Can save significantly on electrical and cooling costs, save rack space
    - It is common to use two, with software RAID 1 for redundancy

**SQL** *intersection*

# Storage Area Networks (SAN)

- **Shared external storage enclosure with multiple components**
  - Large number of drive bays, can usually be expanded
  - Storage processors, large dedicated cache, operating system
  - Usually much higher initial capital cost than DAS
  - Requires some training and expertise to setup and manage
  - Cranky SAN administrator included free of charge!
- **Two main types of SANs**
  - Fiber-channel, using host bus adapter (HBA)
  - iSCSI, using dedicated Ethernet cards
- **SANs are usually optimized for IOPs**
  - Sequential throughput can be severely limited by the interface

SQL
*intersection*

# SAN Administrator Considerations

- **Make an effort to really communicate with your SAN Admin**
  - Let the SAN administrator know the type of workload you have
    - SQL Server OLTP, for example
  - Don't just give the SAN administrator a space requirement!
    - Try to give them useful performance and SLA requirements
- **Your SAN Admin may have different priorities than you**
  - Has to worry about multiple servers with different workloads
  - Has to worry about running low on space in the SAN
  - Has to worry about DBAs complaining about performance
  - Has to worry about CIO complaining about the capital cost

SQL
*intersection*

# SAN Performance Considerations

- **Consider the complete data path to the SAN**
  - HBA/NIC, switches, SAN ports, etc.
  - Be prepared for inconsistent performance with a shared SAN
  - SANs are not magic: the hardware details still matter!
  - SANs are typically sequential throughput limited

# SMB 3.0/3.02 File Shares

- **Server Message Block (SMB) 3.0/3.02**
  - SQL Server 2012+ can store user/system databases on SMB 3.0 file shares
  - SQL Server 2012+ can use SMB 3.0 for traditional FCI instances that require shared storage (without using a SAN)
  - Windows Server 2012 has SMB Direct, which supports the use of network adapters that have Remote Direct Memory Access (RDMA) capability
    - RDMA capable network adapters can function at full speed with very low latency while using very little CPU time on the host
  - Microsoft's Jose Barreto is a great resource about SMB file shares
    - http://blogs.technet.com/b/josebda/

# Negotiated Versions of SMB

| OS | Windows 8.1 WS 2012 R2 | Windows 8 WS 2012 | Windows 7 WS 2008 R2 | Windows Vista WS 2008 | Previous Versions |
|---|---|---|---|---|---|
| Windows 8.1 WS 2012 R2 | **SMB 3.02** | **SMB 3.0** | SMB 2.1 | SMB 2.0 | SMB 1.0 |
| Windows 8 WS 2012 | **SMB 3.0** | **SMB 3.0** | SMB 2.1 | SMB 2.0 | SMB 1.0 |
| Windows 7 WS 2008 R2 | SMB 2.1 | SMB 2.1 | SMB 2.1 | SMB 2.0 | SMB 1.0 |
| Windows Vista WS 2008 | SMB 2.0 | SMB 2.0 | SMB 2.0 | SMB 2.0 | SMB 1.0 |
| Previous Versions | SMB 1.0 | SMB 1.0 | SMB 1.0 | SMB 1.0 | SMB 1.0 |

SQL
*intersection*

# Storage Spaces Direct in Windows Server 2016

- **What is Storage Spaces Direct?**
  - Software-defined storage that is highly available and scalable
  - Storage for Hyper-V and SQL Server
- **Why Storage Spaces Direct ?**
  - Servers with local storage (PCIe, SAS, SATA)
  - Industry standard commodity hardware
  - Lower cost flash storage with SATA SSDs
  - Better flash performance with NVMe SSDs
  - Can use hybrid storage configurations (NVMe, SATA SSD, SATA HDD)
  - Uses Ethernet/RDMA network as storage fabric

**SQL** *intersection*

# Storage Spaces Direct Diagram

# Storage Spaces Direct Hardware Requirements

- **Windows Server certified servers and components**
- **Homogenous server configuration, 4-16 servers in software storage bus**
- **Storage node requirements**
  - Two-socket, Intel Xeon E5-2600 v3 or newer, 128GB of RAM
  - 10GbE or better, RDMA capable NIC strongly recommended
  - Minimum of two caching devices, four capacity devices
  - Simple HBA and expander for SATA/SAS devices
  - RAID controllers are not supported, FC/iSCSI not supported
  - MPIO not supported (single-path SAS is ok)

SQL
*intersection*

# Each Storage Node Has a Built-In Cache

- **Integral part of Software Storage Bus**
- **Cache is coped to local machine, agnostic to storage pools and virtual disks**
- **Read and write cache (depending device type)**
  - Write + read caching when cache device is SSD and capacity device is HDD
  - Write only caching when cache device is SSD and capacity device is SSD
- **Automatic configuration when enabling S2D**
  - Special partition on each caching device
  - Leaves 32GB for pool and virtual disk metadata
  - Round robin binding of SSD to HDD, rebinding with topology change

SQL
*intersection*

# S2D Volume Types

- **Mirror volume**
  - Optimized for performance, all data is hot, least storage space efficiency (33%)
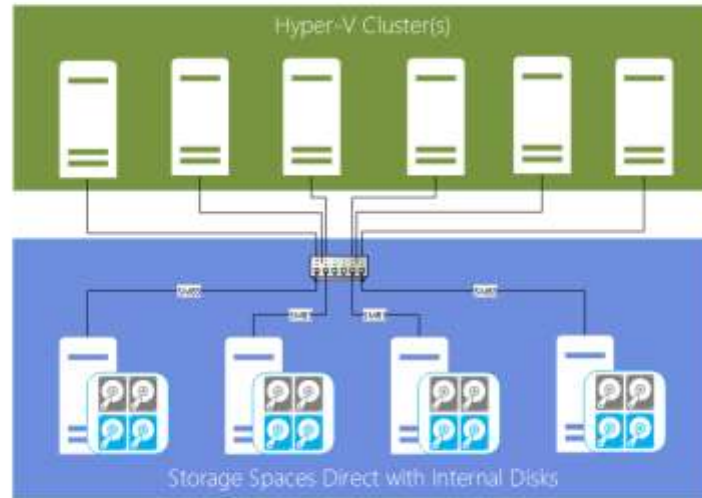  - Uses ReFS or NTFS file system, requires at least two nodes
- **Parity volume**
  - Optimized for capacity, all data is cold, most storage space efficiency (57+%)
  - Uses ReFS or NTFS file system, requires at least four nodes
- **Multi-Resilient volume**
  - Optimized for balance, mix of hot and cold data, efficiency depends on mix
  - Uses ReFS file system, requires at least four nodes

# Disaggregated Deployment with S2D

# Storage Spaces Direct References

- **Storage Spaces Direct in Windows Server 2016 Technical Preview**
  - http://bit.ly/1oc9053
- **Storage Spaces Direct in Technical Preview 4**
  - http://bit.ly/1PvUfQs
- **Hardware options for evaluating Storage Spaces Direct in Technical Preview 4**
  - http://bit.ly/1Mj6YWX
- **Storage Spaces Direct – Under the hood with the Software Storage Bus**
  - http://bit.ly/1PvUxqK

SQL
*intersection*

# Considering Your Workload for Storage

- **SQL Server can have several different common workload types**
  - Online Transaction Processing (OLTP)
  - Reporting against OLTP database(s)
  - Relational Data Warehouse (DW)
  - Online Analytical Processing (OLAP)
- **These workload types have different I/O access patterns**
  - Read/write ratio against different file types
  - Sequential vs. random reads and writes

SQL
*intersection*

# OLTP Workload I/O Access Patterns

- **OLTP workload has frequent writes to data files and log file**
  - Also has frequent reads from data files if the database does not fit in memory
  - Random I/O performance is very important
- **Writes to a single database log file are sequential**
  - Once you have multiple databases with log files on the same LUN, the write activity becomes more random

# DW and Reporting Workload I/O Access Patterns

- **DW/Reporting workload has large sequential reads from data files**
  - Frequent reads from data files if the database does not fit in memory
  - Very little use of log file (except during data loads)
  - Sequential read I/O performance is very important

SQL
*intersection*

# OLAP Workload I/O Access Patterns

- **OLAP workload has lots of random reads from cube files**
  - Random I/O performance is very important
  - Sequential write performance to cube files is important during cube generation

SQL
*intersection*

# RAID Levels and SQL Server Workloads

- **Consider your SQL Server workload type(s)**
  - It directly affects your desired RAID level
  - RAID 10 is better for write-intensive workloads
- **Different types of workloads have different I/O patterns**
  - Percentage of reads/writes, sequential vs. random I/O
  - Use DMV query metrics to determine this, don't guess
- **Different SQL Server file types have different I/O patterns**
  - Data files, log files, tempdb files, backup files, etc.
- **You also need to consider your availability requirements**
  - Some RAID levels are more robust than others: RAID 10 > RAID 50 > RAID 5

SQL
*intersection*

# Selecting a RAID Level For Your SLA

- **RAID is not a substitute for a good backup/restore plan!**
  - No matter what anyone in your organization tells you…
- **RAID is not a substitute for an effective HA/DR strategy**
  - No matter what any vendor tells you…
- **An appropriate RAID level reduces the chance of unplanned downtime**
  - It also reduces the chance of data loss due to disk failure(s)
- **RAID 10 and 50 are the most robust common RAID levels**
  - RAID 5 can only lose one disk in an array before the array is lost
  - Having a higher number of disks in a RAID 5 array increases the statistical chances that any one disk will fail

**SQL** *intersection*

# Choosing Storage Types Based on Workload Type

- **Flash-based storage gives great random I/O performance**
  - It also gives better sequential performance than magnetic storage
  - Flash-based storage is the most expensive storage (per GB), but prices are declining rapidly and approaching parity with magnetic storage
- **Magnetic storage gives fair sequential performance**
  - Magnetic storage gives quite poor random I/O read and write performance
  - Large controller caches can help mask poor random I/O write performance
- **Flash-based storage is the best choice if you have the budget**
  - Use flash-based storage where you have heavy random I/O
  - Use flash-based storage where you have any type of I/O bottlenecks

**SQL**
*intersection*

# Configuring Storage for SQL Server File Types

- **SQL Server data files**
  - Common to use magnetic storage (flash becoming more popular as cost declines)
  - Most common to use RAID 5, 50, or 10
- **SQL Server log files**
  - Common to use magnetic storage (flash becoming more popular as cost declines)
  - Most common to use RAID 10
- **SQL Server tempdb data and log files**
  - Common to use flash-based storage
  - Most common to use RAID 10
  - Acceptable to use magnetic storage if your workload does not heavily use tempdb

**SQL**
*intersection*

# HA/DR Effects on Storage Choices

- **Traditional FCI requires some form of shared storage**
  - Usually a SAN, but SMB 3.0/3.02 can be used with SQL Server 2012 or newer
  - SQL Server 2012 or newer can use local storage for tempdb files with FCI
    - Often a good use for flash storage
    - Better performance and reduces load on the SAN
- **AlwaysOn AGs must use the Windows clustering feature**
  - Can use shared storage, such as a SAN or SMB 3.0/3.02, but is not required
  - Can also use any type of non-shared storage
- **Other HA/DR technologies can use any type of storage**
  - Consider using non-shared storage to eliminate the single point of failure, and get more consistent performance

**SQL**
*intersection*

# Sizing Your Storage Subsystem

- **Use a RAID calculator to ensure you have more than enough disk space**
  - Consider performance advantages of "short-stroking" for magnetic storage
  - Flash-based storage also benefits from ample free space
- **After you have enough space, concentrate on performance**
  - Don't negotiate with yourself! Ask for flash-based storage, ask for RAID 10
  - Consider your workload as you make budget-driven compromises
- **Aim for 10,000-20,000 or more IOPS on all LUNs**
  - More is always better
- **Aim for 1GB/sec or more of sequential throughput on all LUNs**
  - This gives you good performance for common administrative tasks

# Solid State Drives (SSDs)

- **SSD access time does not depend on moving parts**
  - Access time is very fast and consistent across cells
  - Excellent for random I/O reads and writes
- **PCIe flash AIC storage cards allow for much higher sequential throughput**
  - Bypasses traditional SAS/SATA interface bandwidth limitations
- **SSDs are enterprise ready for SQL Server usage**
  - We have many enterprise clients running on them
  - Don't just put tempdb or transaction logs on SSDs!
  - Don't ignore index fragmentation when using SSDs!

**SQL**
*intersection*

# Flash Storage Interfaces and Protocols

- **The interface and protocol have a huge affect on flash performance**
  - SATA/SAS
    - 3Gbps, 6Gbps, 12Gbps (275MB/sec, 550MB/sec, and 1100MB/sec)
    - Typically use Advanced Host Controller Interface (AHCI) protocol
    - AHCI limited to one command queue, 32 commands per queue
  - PCI Express (PCIe)
    - Uses PCIe slot on mother board (or in a front drive bay on some new servers)
  - PCIe Non-Volatile Memory Express (NVMe) protocol
    - NVMe protocol has 65535 queues, 65536 commands per queue
    - Much lower latency and CPU utilization than AHCI

**SQL** *intersection*

# Flash Storage NAND Types

- **Different NAND types affect latency, endurance and cost**
  - Single-Level Cell (SLC)
    - Lowest latency, highest endurance, highest cost
  - Multi-Level Cell (MLC)
    - Higher latency than SLC, lower endurance, lower cost
  - Triple-Level Cell (TLC)
    - Higher latency than MLC, lower endurance, lower cost
    - Write latency suffers more with MLC and TLC
  - Common to use SLC NAND for a small cache in front of MLC or TLC NAND

SQL
*intersection*

# Magnetic Storage vs. Flash-Based Storage

- **Magnetic storage has fair sequential performance**
  - 100-200MB/sec per disk
- **Magnetic storage has very poor random I/O performance**
  - 100-200 IOPS per disk
- **Flash-based storage has very good sequential performance**
  - 12Gbps SAS/SATA does 1100MB/sec, 6Gbps SAS/SATA does 550MB/sec
  - PCIe storage cards can do up to about 6.5GB/sec per card
- **Flash-based storage has excellent random I/O performance**
  - 6Gbps SAS/SATA drives can do about 100,000 IOPS
  - PCIe AIC flash storage cards can do up to about 1.3 million IOPS

**SQL** *intersection*

# Typical HDD Metrics vs. SSD Metrics

| Metric | 15K  Hard Drive | SATA 3 SSD |
|---|---|---|
| Capacity | 900GB | 400GB |
| Average Read Latency (us) | 2000 | 50 |
| Read Bandwidth (MB/s) | 202 MB/s | 550 MB/s |
| Read IOPS (4K QD32) | 150-200 | **90,000** |
| Power (Active/Idle) | 4.25W | 5.2W/0.6W |

**SQL** *intersection*

# Comparative Sequential QD32 Performance

| Drive Type | Sequential Reads | Sequential Writes |
|---|---|---|
| (2) 15K Magnetic SAS in RAID 1 | 154.6 MB/s | 126.2 MB/s |
| (6) 15K Magnetic SAS in RAID 10 | 531.7 MB/s | 414.6 MB/s |
| 1TB Samsung 850 EVO SATA 3 SSD | 555.8 MB/s | 530.5 MB/s |
| 400GB Intel 750 PCIe NVMe AIC | 2382.3 MB/s | 1076.2 MB/s |
| 512GB Samsung 950 PRO M.2 NVMe card | 2598.4 MB/s | 1524.7 MB/s |
| 640GB Fusion-io Duo MLC PCIe AIC | 746.6 MB/s | 516.2 MB/s |
| 2.41TB Fusion-io ioDrive2 Duo PCIe AIC | 1488.6 MB/s | 966.5 MB/s |
| 1.6TB LSI Nytro WarpDrive PCIe AIC | 1483.1 MB/s | 1520.2 MB/s |

**SQL**
*intersection*

# Comparative Random 4K QD32 Performance

| Drive Type | Random Reads | Random Writes |
|---|---|---|
| (2) 15K Magnetic SAS in RAID 1 | 790 IOPS | 1028 IOPS |
| (6) 15K Magnetic SAS in RAID 10 | 2867 IOPS | 3048 IOPS |
| 1TB Samsung 850 EVO SATA 3 SSD | 68,908 IOPS | 70,954 IOPS |
| 400GB Intel 750 PCIe NVMe AIC | 199,216 IOPS | 175,012 IOPS |
| 512GB Samsung 950 PRO M.2 NVMe card | 189,250 IOPS | 103,901 IOPS |
| 640GB Fusion-io Duo MLC PCIe AIC | 49,678 IOPS | 52,655 IOPS |
| 2.41TB Fusion-io ioDrive2 Duo PCIe AIC | 85,815 IOPS | 105,365 IOPS |
| 1.6TB LSI Nytro WarpDrive PCIe AIC | 98,259 IOPS | 78,641 IOPS |

**SQL** *intersection*

# Improving I/O Performance

- **Improving I/O performance at multiple levels**
  - Server hardware evaluation and selection
  - Server hardware configuration settings
  - Storage hardware evaluation and selection
  - Storage hardware configuration settings
  - Operating system configuration settings
  - SQL Server instance-level configuration settings
  - SQL Server database property settings
  - Index tuning
  - Workload tuning

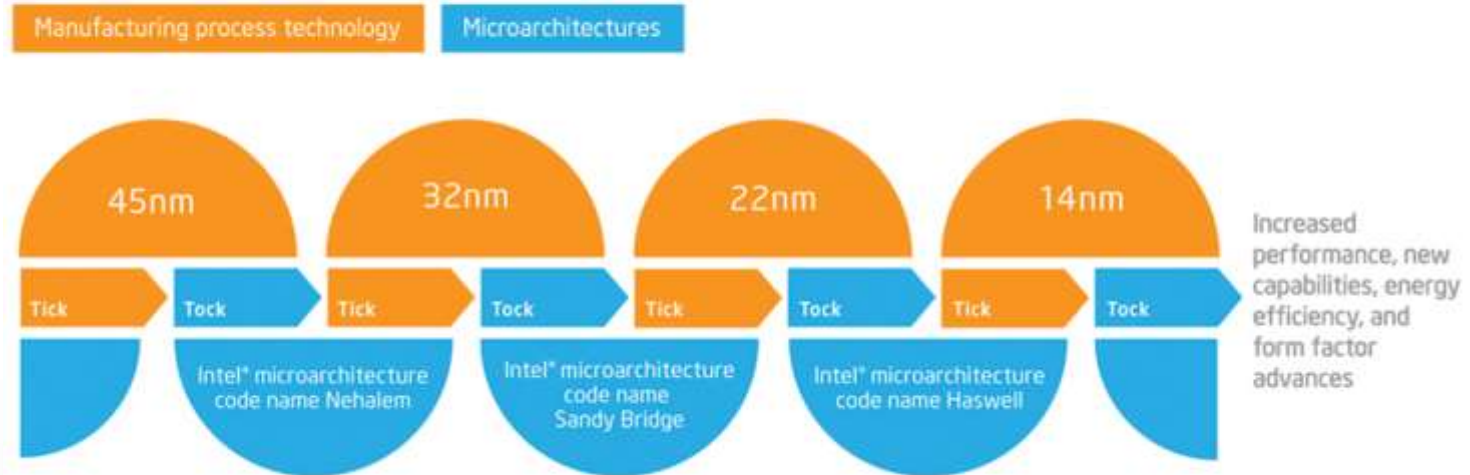# Server Hardware Evaluation and Selection

- **Very important to have relatively modern server hardware**
  - Intel-based servers have many performance and scalability advantages
- **Try to have Xeon E5 or newer (Sandy Bridge-EP or newer)**
  - PCIe 3.0 support, good memory capacity and performance
  - Up to 768GB of RAM with 32GB DIMMs in a two-socket server
  - Up to 9.6 GT/sec QPI speed
- **Try to have Xeon E7 v2 or newer (Ivy Bridge-EX or newer)**
  - PCIe 3.0 support, good memory capacity and performance
  - Up to 3TB of RAM with 32GB DIMMs in a four-socket server
  - Up to 9.6 GT/sec QPI speed

**SQL** *intersection*

# General Processor Considerations

- **Purposely over-provision processors (if you have the budget)**
  - Better single-threaded performance is very important
  - Higher core counts increase overall capacity and scalability, but increase SQL Server license costs
  - Consider lower core count, "frequency optimized" processor models
- **Processors are relatively inexpensive**
  - Adding I/O performance capacity is often more expensive than a good processor
  - The license cost per core is the same, so pick the right processor!
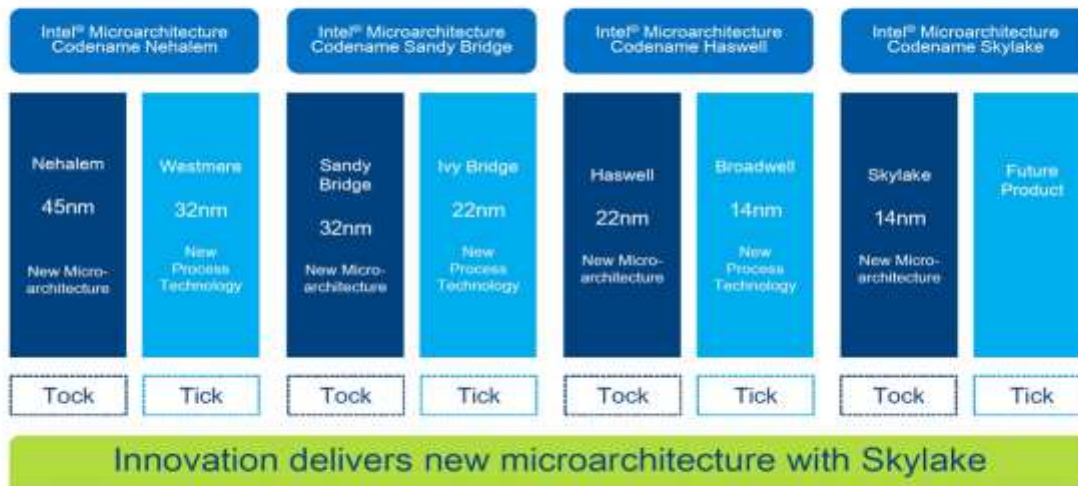  - Don't pick a lower clock speed CPU with the same core count to save money

SQL
*intersection*

# Intel Tick Tock Release Strategy

# Current Intel Tick Tock Release Strategy

# Intel Broadwell-EP Family (Xeon E5-2600 v4)

- **Intel Tick release (Q1 2016)**
  - 14nm process, up to twenty-two cores, up to 55MB L3 cache
  - PCIe 3.0, QPI 1.1, four DDR4 memory controllers
- **Replacement for Haswell-EP(Xeon E5-2600 v3 series)**
  - More physical cores than Ivy Bridge-EP
  - Also has lower core count, frequency-optimized models
- **Xeon E5-2699 v4 is the "top of the line" model**
  - 2.2GHz base, Turbo Boost 2.0 to 3.6GHz, 55MB L3 cache
  - Twenty-two cores, plus hyper-threading, 145W TDP
  - DDR4 2400 support, twelve DIMMs per socket

SQL
*intersection*

# Preferred Broadwell-EP Processors – High Core Count

| Model | Cores/L3 Cache | Base Speed | Turbo Speed | Price |
|-------|----------------|------------|-------------|-------|
| E5-2699 v4 | 22/55 MB | 2.2 GHz | 3.6 GHz | $4,115.00 |
| E5-2698 v4 | 20/50 MB | 2.2 GHz | 3.6 GHz | $3,226.00 |
| E5-2697 v4 | 18/45 MB | 2.3 GHz | 3.6 GHz | $2,702.00 |
| E5-2697A v4 | 16/40 MB | 2.6 GHz | 3.6 GHz | $2,891.00 |
| E5-2690 v4 | 14/35 MB | 2.6 GHz | 3.6 GHz | $2,090.00 |

SQL
*intersection*

# Preferred Broadwell-EP Processors – Low Core Count

| Model | Cores/L3 Cache | Base Speed | Turbo Speed | Price |
|-------|----------------|------------|-------------|-------|
| E5-2687W v4 | 12/30 MB | 3.0 GHz | 3.5 GHz | $2,141.00 |
| E5-2640 v4 | 10/25 MB | 2.4 GHz | 3.4 GHz | $939.00 |
| E5-2667 v4 | 8/25 MB | 3.2 GHz | 3.6 GHz | $2,057.00 |
| E5-2643 v4 | 6/20 MB | 3.4 GHz | 3.7 GHz | $1,552.00 |
| E5-2637 v4 | 4/15 MB | 3.5 GHz | 3.7 GHz | $996.00 |

SQL
*intersection*

# Intel Haswell-EP Family (Xeon E5-2600 v3)

- **Intel Tock release (Q3 2014)**
  - 22nm process, up to eighteen cores, up to 45MB L3 cache
  - PCIe 3.0, QPI 1.1, four DDR4 memory controllers
- **Replacement for Ivy Bridge-EP (Xeon E5-2600 v2 series)**
  - 50% more physical cores than Ivy Bridge-EP
  - Also has lower core count, frequency-optimized models
- **Xeon E5-2699 v3 is the "top of the line" model**
  - 2.3GHz base, Turbo Boost 2.0 to 3.6GHz, 45MB L3 cache
  - Eighteen cores, plus hyper-threading, 145W TDP
  - DDR4 2133 support, twelve DIMMs per socket

SQL
*intersection*

# Preferred Haswell-EP Processors

| Model | Cores/L3 Cache | Base Speed | Turbo Speed | Price |
|-------|----------------|------------|-------------|-------|
| E5-2699 v3 | 18/45 MB | 2.3 GHz | 3.6 GHz | N/A |
| E5-2698 v3 | 16/40 MB | 2.3 GHz | 3.6 GHz | N/A |
| E5-2697 v3 | 14/35 MB | 2.6 GHz | 3.6 GHz | $2,702.00 |
| E5-2690 v3 | 12/30 MB | 2.6 GHz | 3.5 GHz | $2,094.00 |
| E5-2660 v3 | 10/15 MB | 2.6 GHz | 3.3 GHz | $1,449.00 |
| E5-2667 v3 | 8/20 MB | 3.2 GHz | 3.6 GHz | $2,057.00 |
| E5-2643 v3 | 6/20 MB | 3.7 GHz | 3.7 GHz | $1,552.00 |
| E5-2637 v3 | 4/15 MB | 3.5 GHz | 3.7 GHz | $996.00 |

**SQL**
*intersection*

# Intel Haswell-EX Family (Xeon E7-4800/8800 v3)

- **Intel Tock release (May 2015)**
  - 22nm process, up to eighteen cores, up to 45MB L3 cache
  - PCIe 3.0, QPI 1.1, four DDR4 memory controllers
- **Replacement for Ivy Bridge-EX (Xeon E7-4800/8800 v2 series)**
  - 20% more physical cores than Ivy Bridge-EX
  - Intel Transactional Synchronization New Instructions (TSX-NI)
- **Xeon E5-8890 v3 is the "top of the line" model**
  - 2.5GHz base, Turbo Boost 2.0 to 3.3GHz, 45MB L3 cache
  - Eighteen cores, plus hyper-threading, 165W TDP
  - DDR4 1833 support, 24 DIMMs per socket

SQL *intersection*

# Preferred Haswell-EX Processors

| Model | Cores/L3 Cache | Base Speed | Turbo Speed | Price |
|-------|----------------|------------|-------------|-------|
| E7-8890 v3 | 18/45 MB | 2.5 GHz | 3.3 GHz | $7,175.00 |
| E7-8867 v3 | 16/45 MB | 2.5 GHz | 3.3 GHz | $4,672.00 |
| E7-4850 v3 | 14/35 MB | 2.2 GHz | 2.8 GHz | $3,003.00 |
| E7-4830 v3 | 12/30 MB | 2.1 GHz | 2.7 GHz | $2,170.00 |
| E7-8891 v3 | 10/45 MB | 2.8 GHz | 3.5 GHz | $6,841.00 |
| E7-8893 v3 | 4/45 MB | 3.2 GHz | 3.5 GHz | $6,841.00 |

SQL
*intersection*

# General Memory Considerations

- **Maximize your physical RAM (within SQL Server license limits)**
  - Larger buffer pool cache reduces physical reads from disk subsystem
    - More data in the buffer pool (logical vs. physical reads)
  - RAM is faster than any disk subsystem
    - RAM may be less expensive than enterprise-class storage
    - Orders of magnitude difference in latency
  - Can reduce the frequency of lazy writes and checkpoints
    - Helps even out the write workload to your data files

# DDR4 PC4-17000 ECC Memory Prices

- **32GB module**          **$250.00**                    **$7.81/GB**
- **16GB module**          **$115.00**                    **$7.19/GB**
- **8GB module**           **$  62.00**                   **$7.75/GB**
  - Retail prices from Crucial.com  (3/31/2016)
- **Current capacity/price sweet spot is 32GB modules!**

# Server Hardware Configuration Settings

- **Intel Hyper-threading**
  - Gives 20-30% more CPU capacity, not more performance
  - Should be enabled with most workloads
    - Exceptions: Virtualization host, some DW or reporting workloads
- **Intel Turbo Boost**
  - Boosts clock speed of individual cores, should always be enabled
  - Make sure to check processor core speed with CPU-Z
- **BIOS or UEFI power management settings**
  - Set to OS control or disable. BIOS-level power management will override the Windows Power Plan setting

**SQL**
*intersection*

# CPU-Z

**Demo**

Using CPU-Z

SQL
*intersection*

# Storage Hardware Evaluation and Selection

- **Do your research, and check for online reviews from established sites**
  - StorageReview
  - SSDReview
- **Run your own storage benchmarks on an evaluation unit if possible**
  - Definitely run storage benchmarks on each logical drive before you install SQL Server
- **Do more storage testing after you install SQL Server**
  - Using a representative workload if possible
  - Doing common administrative tasks

**SQL**
*intersection*

# Storage Hardware Configuration Settings

- **Consult any available documentation from the storage vendor**
  - Many vendors have SQL Server specific guidance for their product
  - Several SQL Server "personalities" work for major storage vendors
    - Jimmy May – SanDisk
    - Argenis Fernandez – Pure Storage

# Operating System Configuration Settings

- **Change Windows Power Plan to High Performance**
  - Default setting is still Balanced, even in Windows Server 2016
  - Can have a substantial effect on processor and storage performance
- **Enable Windows Instant File Initialization**
  - Huge effect on database restore times, database data file growth times
  - Grant "Perform volume maintenance tasks" right to SQL Server Service account
- **Use Lock Pages in Memory (LPIM)**
  - Prevents OS from paging out SQL Server data
  - Grant "Lock pages in memory" right to SQL Server Service account

# SQL Server Instance-Level Configuration Settings

- **Several sp_configure settings should be changed from default values**
  - Backup checksum default should be 1 (new in SQL Server 2014)
  - Backup compression default should be 1 (in most cases)
  - Cost threshold for parallelism often should be raised above 5 (in most cases)
  - Max degree of parallelism should be set to the number of cores in a NUMA node
  - Max server memory should be set to an appropriate value based on the workload
  - Optimize for ad hoc workloads should be 1

SQL
*intersection*

# Demo

**Changing instance-level configuration settings**

SQL *intersection*

# SQL Server tempdb Settings

- **Start with 4-8 tempdb data files, to reduce allocation contention**
  - Make sure they are all the same size, with the same auto growth setting
  - Make sure they are not using percent growth
  - In most cases, it is ok for them to be on the same LUN, and to have the tempdb log file in the same location
- **Use dedicated, fast local storage if possible**
  - Consider using flash storage if your workload requires it
- **Make sure to enable TF 1118**
  - Not necessary in SQL Server 2016

**SQL** *intersection*

# SQL Server Database Property Settings (1)

- **Use a MAIN file group with at least two data files with new databases**
  - Make the MAIN file group the default file group. This will put your user objects in the MAIN file group and the system objects in the PRIMARY file group
  - This gives you more flexibility to lay out your data files and potentially more I/O performance
- **Make sure to use a reasonable auto growth size in MB**
  - Try to manually manage file growth, but leave autogrow enabled
  - Don't use percent growth for data or log files

**SQL**
*intersection*

# SQL Server Database Property Settings (2)

- **Control the VLF counts on your database log file**
  - Manually grow the log file in larger chunks (1000, 2000, or 4000MB)
  - Set auto grow to a similar large growth increment
- **Keep VLF counts below 100-200 (depending on log file size)**
  - Full recovery model
    - Transaction log backup, then shrink the log file. May have to repeat several times
    - May have to generate some log activity to get the log file to shrink
  - Simple recovery model
    - Checkpoint, then shrink the log file. May have to repeat several times

# SQL Server Database Property Settings (3)

- **Auto update statistics asynchronously**
  - Should be enabled with most workloads
  - Make sure you have a new enough build of SQL Server for your major version
- **Delayed durability**
  - New for SQL Server 2014
  - Can noticeably improve transaction rates if writing to transaction log is your main bottleneck
  - Some risk of data loss, so understand how it works!

SQL
*intersection*

# ALTER DATABASE SCOPED CONFIGURATION

- **New feature in SQL Server 2016, lets you control database-level configuration**
  - Enable/disable legacy cardinality estimation, independent of compatibility level
  - Enable/disable parameter sniffing
  - Enable/disable query optimizer hotfixes (TF 4199)
  - Set max degree of parallelism for individual databases
  - The four options above can be set for the primary replica or for secondary replicas
  - Clear the plan cache for an individual database
- Using ALTER DATABASE SCOPED CONFIGURATION in SQL Server 2016
  - http://bit.ly/1Rvok7b

SQL *intersection*

# Demo

**Using ALTER DATABASE SCOPED CONFIGURATION**

**SQL**
*intersection*

# Index Tuning

- **Proper index tuning can have huge positive performance benefits**
  - Reduced CPU, memory, and I/O pressure
  - Reduced query execution times
  - Reduced locking, blocking and deadlock issues
- **Consider your overall workload and individual table volatility**
  - Indexes help read performance, but hurt write performance
  - Look for missing indexes and missing index warnings in the plan cache
  - Look for unused or duplicate indexes
- **Consider data compression or clustered columnstore indexes**

SQL
*intersection*

# Demo

## Index Tuning Queries

**SQL** *intersection*

# Demo

**Using SQL Server Data Compression**

SQL *intersection*

# Demo

## Using SQL Server Clustered Columnstore Indexes

**SQL** *intersection*

# Workload Tuning

- **Use my "bad man list" DMV queries to find expensive stored procedures**
  - SP execution counts, SP avg elapsed time, SP worker time, SP logical reads
- **Focus on the top five stored procedures in each list**
  - Prioritize the area where instance is under the most stress
  - Make it a team effort and an iterative process, and repeat as necessary

SQL
*intersection*

# Demo

**Bad Man List queries**

SQL *intersection*

# Demo

## Using Resource Governor to limit I/O


SQL *intersection*

# References (1)

- **Windows Server 2012 R2: Which version of the SMB protocol (SMB 1.0, SMB 2.0, SMB 2.1, SMB 3.0 or SMB 3.02) are you using?**
  - http://bit.ly/18uOEI4
- **Updated Links on Windows Server 2012 R2 File Server and SMB 3.02**
  - http://bit.ly/1iJKMWb
- **Storage Review SQL Server OLTP Benchmark**
  - http://bit.ly/1jEDu9m
- **My Pluralsight courses**
  - http://bit.ly/1EUh7v9

**SQL** *intersection*

# References (2)

- **Geekbench**
  - http://bit.ly/UGrGbu
- **TPC-E OLTP benchmark**
  - http://bit.ly/UGs2Pm
- **CPU-Z tool**
  - http://bit.ly/korH23
- **Intel Ark database**
  - http://ark.intel.com/

**SQL** *intersection*

# Review

- **SQL Server has five primary storage types**
  - Internal, PCIe flash, DAS, SANs, and SMB 3.0 file shares with SOFS and S2D
- **Different types of SQL Server workloads affect I/O patterns**
  - OLTP, DW, OLAP, mixed, database maintenance, etc.
- **Different SQL Server file types have different I/O patterns**
  - Data files, log files, tempdb files, backup files, etc.
- **Choose an appropriate RAID level for your workload**
  - You also need to consider your SLA requirements
- **Make sure to consider your sequential throughput**
  - Very important for day-to-day operations and DR requirements

## SQL
*intersection*