

## SQL Server I/O Reliability Program Review Requirements

### Applies To

- SQL Server 2000 (all editions)
- SQL Server 2005 (all editions)
- SQL Server 2008 (all editions)
- SQL Server 2012 (all editions)
- SQL Server 2014 (all editions)
- SQL Server 2016 (all editions)

### Overview

This page outlines the required and recommended behaviors an I/O subsystem must provide for Microsoft SQL Server. [KB 967576](#) provides extended details and reference links related to SQL Server I/O requirements.

**Copyright**

©2015 Microsoft Corporation. All rights reserved. This document is provided "as-is." Information and views expressed in this document, including URL and other Internet Web site references, may change without notice. You bear the risk of using it.

This document does not provide you with any legal rights to any intellectual property in any Microsoft product. You may copy and use this document for your internal, reference purposes.

## Requirement Types

The storage system capabilities are divided into two categories - required or recommended - to provide a guide to the minimum requirement and best practice recommendations.

Definitions are listed below:

- **Required** - A capability or property that the subsystem must provide to pass the requirements of the SQL Server I/O Reliability Review Program.
  - **Recommended** - A capability or property that the subsystem should provide for optimal compatibility with SQL Server.
- ❖ A SQL Server I/O Reliability white paper must document specific product configurations that support each core requirement.
  - ❖ If one or more of the core requirements is not met engage the SQL Server team to discuss requirements for qualification.
  - ❖ If advanced feature requirements are met by the I/O solution those features should be clearly documented.

## Core 1.00: Windows Logo Certification - (Required)

Microsoft Windows logo certification helps ensure the safety of Microsoft SQL Server data by testing various aspects. To be compliant with the SQL Server I/O Reliability Program Review Program, solutions must pass and maintain the latest certifications for Windows logos.

The Windows logo program has both hardware and a software tracks. SQL Server I/O Reliability requires completion and logo certification for the tests applicable to the I/O solution.

## Core 1.01: Core Windows API Support - (Required)

SQL Server utilizes several APIs to enable secure data storage. A storage solution must ensure that a system supports specific API properties throughout the various layers and implementations of the I/O solution. The following table outlines the required storage APIs and outlines the fundamental properties of each API which must be supported.

Core API	Fundamentals
<a href="#">CreateFile</a>	<p>SQL Server uses many of the options allowed by CreateFile for various operations on database, backup, security and other files. The following are specific options of importance.</p> <ul style="list-style-type: none"> <li>• FILE_FLAG_WRITE_THROUGH</li> <li>• FILE_FLAG_NO_BUFFERING</li> <li>• FILE_FLAG_OVERLAPPED</li> <li>• FILE_FLAG_RANDOM_ACCESS</li> <li>• FILE_FLAG_SEQUENTIAL_SCAN</li> <li>• FILE_SHARE_READ</li> <li>• FILE_SHARE_WRITE</li> </ul>

	<ul style="list-style-type: none"> <li>• GENERIC_READ</li> <li>• GENERIC_WRITE</li> <li>• READ_CONTROL</li> <li>• WRITE_DAC</li> <li>• WRITE_OWNER</li> <li>• SECURITY_ANONYMOUS</li> <li>• SECURITY_SQOS_PRESENT</li> <li>• SECURITY_VALID_SQOS_FLAGS</li> <li>• SERVER_SQOS_MANDATORY</li> </ul> <p><b>Recommended:</b> Support for secondary stream usage. Secondary stream usage and naming is outlined in SQL Server I/O Basics – Chapter 2’. (<b>Note:</b> Applies to all SQL Server versions.)</p> <p><a href="#">SQL Server 2014 extends online DBCC behavior for ReFS.</a></p>
<p><a href="#">DeviceloControl</a></p>	<p>DeviceloControl is used to establish and query file system behaviors. Discovering details about the volume information, establishing sparse file settings; obtaining sparse file allocation information or other attributes.</p> <p>The following are specific options and data structures of importance. The I/O system should ensure accuracy of the returned data.</p> <ul style="list-style-type: none"> <li>• FSCTL_SET_SPARSE</li> <li>• FSCTL_SET_ZERO_DATA</li> <li>• FSCTL_QUERY_ALLOCATED_RANGES</li> <li>• FSCTL_SET_COMPRESSION</li> <li>• IOCTL_DISK_GET_LENGTH_INFO</li> <li>• IOCTL_DISK_GET_PARTITION_INFO_EX</li> <li>• IOCTL_DISK_GET_DRIVE_GEOMETRY</li> <li>• GPT_BASIC_DATA_ATTRIBUTE_SHADOW_COPY</li> <li>• GPT_BASIC_DATA_ATTRIBUTE_READ_ONLY</li> </ul> <p><b>Note:</b> DeviceloControl can be called by using and without the OVERLAPPED structure.</p>
<p><a href="#">FlushFileBuffers</a>  <a href="#">GetVolumePathName</a>  <a href="#">GetVolumeInformation</a>  <a href="#">GetVolumeNameForVolumeMountPoint</a></p>	<p>The solution must honor the SL_WRITE_THROUGH and IRP_MJ_FLUSH_BUFFERS.</p> <ul style="list-style-type: none"> <li>• Used to harden a backup before the transaction log is truncated</li> <li>• Used when file size changes are made to ensure metadata is flushed</li> <li>• Enables support of mount points.</li> </ul>

	<ul style="list-style-type: none"> <li>Provides details about the volume, similar to those outlined in the DeviceIoControl entry.</li> </ul>
--	--

Core API	Behavioral Needs
<a href="#">WriteFile</a>	<ul style="list-style-type: none"> <li>OVERLAPPED Present</li> <li>OVERLAPPED Not-Present</li> <li>General I/O size ranges from 512-bytes to 4MB</li> </ul>
<a href="#">WriteFileGather</a>	<ul style="list-style-type: none"> <li>OVERLAPPED Present</li> <li>OVERLAPPED Not-Present</li> <li>General I/O size ranges from 8KB to 8MB</li> <li>Systems that do not support scatter/gather can experience performance degradation</li> <li>Used heavily by SQL Server lazy writer, eager writes (bulk operations) and checkpoint processing</li> <li>WOW64 on X64 support needed</li> </ul>
<a href="#">ReadFile</a>	<ul style="list-style-type: none"> <li>OVERLAPPED Present</li> <li>OVERLAPPED Not-Present</li> <li>General I/O size ranges from 512-bytes to 4MB</li> </ul>
<a href="#">ReadFileScatter</a>	<ul style="list-style-type: none"> <li>OVERLAPPED Present</li> <li>OVERLAPPED Not-Present</li> <li>General I/O size ranges from 8KB to 10MB</li> <li>Systems that do not support scatter/gather can experience performance degradation</li> <li>Used heavily by SQL Server read ahead logic</li> <li>WOW64 on X64 Support needed</li> </ul>

SQLIOSim.exe, a replacement for SQLIOStress.exe, is shipped with SQL Server as a utility for use by storage solution providers to improve testing. A storage solution provider may use SQLIOSim.exe to help review and test an implementation. (Reference Core 1.08 for testing requirements.)

The following is an example of the usage of *FlushFileBuffers*.

The file-system (NTFS) has two ways to tell the subsystem to flush the data to stable media. Marking each write IRP as SL\_WRITE\_THROUGH, and sending IRP\_MJ\_FLUSH\_BUFFERS.

When *FlushFileBuffers* is called, NTFS will write each page remaining in file system cache and invoke IRP\_MJ\_FLUSH\_BUFFERS. The flush must ensure all data has been written to stable media.

Just using FLAG\_WRITE\_THROUGH + FILE\_FLAG\_NO\_BUFFERING does not cause IRP\_MJ\_FLUSH\_BUFFERS, although SL\_WRITE\_THROUGH is set on every IRP.

Certain IDE drives will reorder your writes and making SQL Server susceptible to data loss after a power outage. FlushFileBuffers is sufficient to make sure that data is persisted on disk.

### **Core 1.02: Stable Media - (Required)**

SQL Server relies on the Write-Ahead Logging (WAL) protocol to maintain the Atomicity, Consistency, Isolation, and Durability (ACID) properties of the database. WAL relies on stable media capabilities. A solution must comply with this stable media intention. For detailed information, see the 'Power Outage Testing – Pull The Plug' section in [Microsoft SQL Server I/O Basics Chapter 2](#). (**Note:** Applies to all SQL Server versions.)

### **Core 1.03: Forced Unit Access (FUA) and Write-Through - (Required)**

To support Write-Ahead Logging (WAL), SQL Server uses FILE\_FLAG\_WRITE\_THROUGH when opening database files. SQL Server also uses *FlushFileBuffers* during various operations. Write-through and flushing to stable media must be supported by storage solutions.

All components in a solution must honor the write-to-stable media intent. This includes, but is not limited to, caching components.

It is not enough to honor WAL for SQL Server log files only. Data files and backup streams also depend on WAL behavior.

Many storage products include battery-backed caching mechanisms. If these caching mechanisms are present in the solution, the SQL Server Always On solution white paper should document the practical limits of the write-through stable media protection for a production environment.

For more information, see the links listed in the References section at the end of this paper, and the following Microsoft Knowledge Base article: [KB917043](#) - Key factors to consider when you evaluate third-party file cache systems with SQL Server.

### **Core 1.04: Asynchronous Capabilities - (Required)**

SQL Server performs most of its I/O using asynchronous capabilities. If a request specifies asynchronous operation, no API call should cause a synchronous condition. Synchronous I/O can cause unexpected scheduler and concurrency issues. Therefore, a SQL Server solution must provide asynchronous I/O capabilities.

For more information about how a synchronous action can affect the Microsoft SQL Server scheduler, see the white paper, [How To Diagnosis and Correct Errors 17883, 17884, 17887, and 17888](#).

### **Core 1.05: Write Ordering - (Required)**

A tenant of the WAL protocol is write ordering preservation. Any SQL Server solution must provide write ordering preservation.

For more information about write ordering requirements, see the Write Ordering, FlushFileBuffers, Backup Hardening, and Remote Mirroring sections of the white paper, [SQL Server 2000 I/O Basics](#).

**(Note:** Applies to all SQL Server versions)

### **Core 1.06: Torn I/O Protection - (Required)**

A SQL Server I/O Reliability solution must provide sector alignment and sizing in a way that prevents torn I/O including splitting I/Os across various I/O entities in the I/O path.

Additionally, a solution must accurately report sector size to Windows I/O APIs. Accurately reporting sector size helps prevent sector size mismatches and torn writes. For example, a drive that does 4 KB writes reports 512 bytes while the drive performs a read/write of the 4 KB sectors. This inaccuracy in reporting sector size can create a condition where data is lost or exposed as torn writes. Any SQL Server I/O Reliability solution must document configurations in such a way that use actual sector sizes from the sector size list that is supported by Microsoft SQL Server: 512, 1024, 2048, and 4096 bytes.

To indicate when a torn-write situation occurs, we recommended that the solution log appropriate warning events.

The SQL Server I/O Reliability solution white paper must include information about the configuration requirements needed for the solution to meet the torn I/O requirements.

For more information, see the Torn I/O, Log Parity, and Sector Size sections located in the white paper, [Microsoft SQL Server I/O Basics Chapter 2](#). **(Note:** Applies to all SQL Server versions.)

### **SSD/Flash (Non-Spinning Media)**

Sector sizes become opaque (often simulated) to the I/O stack for solutions which are not based on spinning media. Solutions with non-spinning media designs must still avoid torn writes. These solutions should document how torn writes are avoided. They should also document how to configure the solution to reduce Read/Modify/Write I/O patterns, reduce wear of the media and optimize I/O performance.

### **Core 1.07: NTFS Support - (Required)**

You must support NTFS capabilities. This includes but not limited to the following:

- Sparse Files
- File Streams
- Encryption
- Compression
- All Security Properties

The solution must support sparse files on NTFS based file systems. Microsoft SQL Server 2005 and newer versions use sparse files in support of DBCC CHECK\* commands and snapshot databases.

Common copy and compression utilities may not honor sparse file metadata but instead copy all bytes, ignoring the sparse allocations and requiring full storage space. Storage solution providers may choose to provide utilities to copy or move sparse files without destroying the sparse file intent.

## Core 1.08: Testing - (Required)

The SQL Server I/O Reliability Review Program requires successful execution of the following data durability test suites.

1. The latest SQLIOSim.exe (*installed in BINN directory during SQL Server installation*) for data durability and integrity testing. SQLIOSim.exe is shipped with SQL Server 2008 and newer versions.

The following configuration file is to be used to complete a continuous, 24 hour stress test. To open double click on the file name below.



Alwayson.sqliosim.cfg.ini

2. Common benchmark suites such as (TPC-E, TPC-C, TPC-H) should be executed with:
  - a. Database protection level (PAGE AUDIT) set to checksum.
  - b. Trace flags (818, 815) are enabled

**Note:** BIOS and UEFI bugs can lead to unexpected behavior for the solution. Solution testing must include full power outage restarts followed by DBCC CHECKDB validation.

It is recommended the following tests and reviews be performed.

- x86 with 3GB Enabled (optional)
- x86 with PAE Enabled (optional)
- WOW64 running x86 in x64 (optional)
- Low paged and non-paged pool conditions
- Excessive outstanding request boundaries. For example 10000+.
- Memory requests are not forced to a specific memory location. For example, a 64 bit driver requiring a memory allocation under the 4GB can force aggressive working set trimming.

## Advanced 2.01: Write Ordering - (Required - Remote Storage Solutions)

For remote and mirrored I/O destinations all the paths must honor write ordering across the database files. A SQL Server I/O Reliability solution white paper must include information about the configuration requirements needed for the solution to meet the write, ordering requirement. For example, a solution

that requires a consistency group might specify this configuration requirement as: "All files associated with a database must be configured in a single consistency group."

#### Example Configuration

A solution has the following configuration:

- Data File - Device A - Subsystem #1
- Log File - Device B - Subsystem #2

**Note:** If these subsystems use separate physical paths with different caching, SQL Server may not be able to support this configuration because the caching mechanisms may not present a coherent cache. The subsystems may require a third element to maintain cache coherency across the disparate caches. Consider including system databases in the consistency group to provide enhanced metadata consistency.

The same caching problem described in the example configuration can also occur across network boundaries. If a database backup is written to a UNC path but FlushFileBuffers only ensures that the local system file cache is flushed, SQL Server may be exposed to data loss. Network based solutions must ensure stable media delivery.

#### **Non-Battery Backed Cache Solutions**

If the solution uses cache that is not battery backed it must provide stable media and write ordering guarantees as well. An exception is [TEMPDB](#).

This type of solution often holds the write buffers in RAM memory during transmission and must use a two phase commit approach to maintain stable media and write ordering requirements.

#### **Advanced 2.02: Transactional Sector/Block Rewrites - (Required)**

Solutions involving movement of sectors or blocks must provide transactional safety while maintaining asynchronous capabilities. Sectors or blocks cannot be rewritten or changed unless transactional safety can be guaranteed.

See the Defragmenting Disk Drives and Sector Size sections located in [SQL Server 2000 I/O Basics](#) for complete details. (**Note:** Applies to all SQL Server versions.)

#### **Advanced 2.03: VDI - (Required)**

VDI solutions must meet the requirements outlined in the [SQL Server 2005 Virtual Backup Device Interface \(VDI\) Specification](#).

#### **Advanced 2.04: Clustering - (Required)**

For clustering solutions, the shared disk should be part of the Hardware Compatibility List for Windows Server.

### **Advanced 2.05: File Streams - (Required)**

The SQL Server file streams feature requires NTFS transactional guarantees. Compliance of SQL Server file stream requires the following.

- File system must report NTFS.
- File system must support the NTFS mini-filter stack including support of the filter manager contexts FSRTL\_FLAG2\_SUPPORTS\_FILTER\_CONTEXTS in FCBs.
- Support for Extended Attributes.
- Metadata changes on the same volume must maintain ordering. For example, a metadata change of File A followed by a metadata change of File B will maintain order, even after a crash recovery.
- Directory scan enumeration returns the current state of the directory. For example, if a directory contains files A and C when the scan is started but during the scan file B is added; then A and C should be returned and B is optional. The file system must not look at the count of files at scan startup as an absolute. In the example, the count would be two(2). The scan should not return A and B only. The enumeration behavior of B is undefined but A and C were present at the start of the scan and unaltered during the scan so they should be returned in the scan.
- It is recommended that the solution participates in Microsoft Plug Fest Interop testing.
- File stream access does not support OpLocks.

### **Advanced 2.06: Protection- (Recommended)**

Data durability compromises can frequently be predicted or avoided. This is frequently referred under initiatives such as S.M.A.R.T. Solutions are encouraged to provide advanced data protection features.

[Custom stability checks](#) can be implemented in conjunction with an Always On solution.

### **Advanced 2.07: Hardware Virtualization (Required)**

Solutions involving virtualized environments must comply with Windows [SVVP](#) program as outlined in the [SQL Server virtualization support policy](#).

## SQL Server I/O Patterns 3.00

The following table outlines the common SQL Server I/O patterns. Leverage this table to:

- Evaluate the I/O solution
- Establish appropriate testing patterns
- Document configurations and recommendations

I/O Action	Common I/O Sizes
Database page reads (mdf/ndf)	8K to 512K and may not be block aligned. Accessed in either allocation and index/data order depending on query plan selection Read over write patterns leveraged by SQL Server.
Database page writes (mdf.ndf)	8K to 128K and may not be block aligned Ex: Eager Writes (select into), Checkpoint, Lazy writer
Database initialization writes	8K (may not be block aligned) to 8MB. Zero'ed 8K buffer posted multiple times (WriteFileGather) for multiple offsets.
Log writes	512 (smallest sector size) to 60K, always sector aligned
Log reads	512 (smallest sector size) to 480K
Backup reads/writes	Default 4MB ( <a href="#">based on buffer size calculations</a> )
Restore reads/writes	Default 4MB ( <a href="#">based on buffer size calculations</a> )
In-Memory Database Checkpoint files (reads/writes)	256K, serial access pattern goals and file reuse
In-Memory Database Delta files (reads/writes)	4K, serial access pattern goals
File Stream (T-SQL) access (reads/writes)	8K
File Stream external access (reads/writes)	<i>Application dependent</i>
T-SQL Blob reads	8K to 512K, total outstanding 2048 pages
Column Store (reads/writes)	8K to 16MB, uses T-SQL blob reads, total outstanding 40,000 pages
T-SQL BULK INSERT	64K
BCP.exe writes	<a href="#">Reference details</a>
Sparse file writes	8K
Sparse file reads	See 'database page reads'
Buffer Pool Extension writes	8K to 128K

Buffer Pool Extension reads	8K
-----------------------------	----

**Note:** The SQL Server behavior may vary. The table is intended to outline the most common patterns leveraged.

### **I/O Parallelism**

Numerous SQL Server I/O paths leverage full asynchronous capabilities to keep multiple I/O requests active in parallel. For example, SQL Server Enterprise, database page, read-ahead attempts to keep 5000 page requests outstanding across a total of 78 x 512K requests.