



Microsoft® SQL Server® 2012

SQL Server Technical Article

Writer: Jonathan Foster

Technical Reviewer: Karthik Pinnamaneni; Andre Ciabattari

Published: November, 2013

Applies to: SQL Server 2012

Summary: This paper examines the performance of an XtremIO All Flash array in an I/O intensive BI environment.

Copyright

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.

Some examples depicted herein are provided for illustration only and are fictitious. No real association or connection is intended or should be inferred.

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.

© 2013 Microsoft. All rights reserved.

Introduction

The Customer Service & Support BI group at Microsoft (CSSBI) was offered the opportunity to test BI-related SQL loads on EMC's XtremIO storage unit for ultra-fast IO rates. The group decided to create a load test that simulated peak activity in their data acquisition platform (DAP). The DAP is a single SQL Instance containing 100's of various size databases hosted on a very powerful server. These databases are updated w/ new data at various frequencies while at the same time queried by managed systems, reports, and ad hoc users. This paper presents the disk I/O performance of the XtremIO unit in a DAP test environment and also provides disk I/O performance metrics from a standard SAN to be used as a reference point. The purpose of this paper is to provide BI architects and DBA's with viable performance information for helping select viable hardware configurations to meet specific performance needs.

Background

As mentioned previously, the CSSBI DAP server is subjected to competitive IO requests throughout each day with the heaviest load occurring in the early morning business hours. Many of these IO transfers involve large datasets like clickstream or user event data and/or large-object data types like online forum chats and XML data. A portion of these large data flows are repetitive meaning the same query types occur multiple times each day. Near-real time (NRT) subsets of these data flows are also occurring simultaneously. CSSBI's value to stakeholders is its ability to provide consistently good performance in an easily accessible symmetric multiprocessing (SMP) environment. The project's success has resulted in exponential growth of both data and users which has forced the IT owners to regularly upgrade the underlying infrastructure & platform to ensure optimal performance.

Purpose

When EMC shared the technical details and future integration plans of the XtremIO storage unit, the CSSBI IT team wanted to evaluate whether the product might solve a specific need for expanding the NRT or real-time (RT) capability of the DAP for larger datasets as described previously. Most often in these scenarios the solution is one that involves more expensive massive parallel processing (MPP) capability. The CSSBI IT team's goal is to try and meet the business requirements with a less expensive SMP capability for as long as possible.

Test Environment

Hardware	
Platform	Cisco B230-Base-M2 Blade
CPU (Clock Speed, Cache, Max TDP)	2 x 10 core (hyper threading) Intel Xeon E72860 @ 2.27GHz, 24MB L3 Cache, 130W (40 logical processors)
Memory	28 x 8GB DDR3 (224GB total)
HDDs (Capacity, Rotational Speed, Interface, Form Factor)	Local Storage 2 x 64GB SSD (Raid 1) (C: / D:)
XtremIO Storage	1 XBrick @ 7.5TB usable presented as 5 separate volumes of 1TB (1) & 1.5TB(4)
SAN Storage	EMC VNX 5700 having 1 auto-tiered pool composed of SAS 10k (95%) & SSD (5%)
Controller	Cisco N20-AC0002 Con adapter (dual 8GB FC paths)
Power Supply	(Power supplies in chassis not blades)
Network	Single Dual Port Embedded Cisco NIC @ 10Gbps

Storage Configuration	
XtremIO X-Brick	XtremIO Data Protection methodology (XDP)
SAN Pool	RAID 5

Software	
OS	Microsoft Windows Server 2012 Datacenter
SQL	Microsoft SQL Server 2012 Developer SP1 (X64)

Test

The test involved a variety of DDL & DML operations having distinctively different I/O patterns organized to run predominately synchronously in order to force longer disk queues. The test was designed and executed within SSIS with the beginning and ending of each step logged to audit tables. The relational landscape consisted of 9 databases containing a mixture of clickstream, incident management, and chat data with a storage footprint of approximately 3.5TB. Other important environment details to call out are:

- The source tables all contained clustered indexes having various levels of fragmentation.
- All SQL engine, system and user db files were intentionally placed on volumes hosted by 1 storage device. For example, during the XtremIO testing all SQL and db files were placed on volumes presented by the XtremIO device. Conversely, during the SAN comparison test all SQL and db files were placed on volumes presented by the SAN.
- SQL Server memory was intentionally fixed @ 190GB for min and max settings to eliminate memory allocation swaps.
- SQL system cache & buffer pool were cleansed prior to each single test run
- All data flows used OLE DB connections with Rows Per Batch set to 10k & Maximum Insert Commit Size set to 100k
- The clickstream database totaled 1.7TB and was composed of 3 monthly partitions
- The incident & chat databases which contained the bulk of LOB data totaled 1.9TB w/ each having their largest tables spread in monthly partitions

The following pages summarize each individual test component in a graphical format. The color and pattern of each graphic will correspond to individual line items in the disk latency chart presented later in the results.

Click stream test flow

Remote extract

Remote extract & bulk insert 22.7 million rows totaling 14.6 GB and multicast to 2 heaps (divided into 4 asynchronous data flows)

Local extract

Local extract & bulk insert 79 million rows totaling 22.9 GB (divided into 4 synchronous SQL tasks)

Complex aggregation

1 query scans 17 million rows within a specific range (clustered index seek), inserts 6.5 million rows into a temp table, then derives a summary result set requiring several temp table subsets, joins, groupings, sorts, and index creation

Sequential READs

100 asynchronous queries each using a distinct WHERE filter that performs a clustered index seek on approximately 10 million rows

Random READs

10 million asynchronous queries each using a distinct WHERE filter that performs a clustered index seek on approximately 1 to 10 rows

Simple aggregation

1 query scans 4 varchar columns in a range of 693 million rows into a CTE, then performs a series of parsing operations involving temp tables, joins, groupings, sorts

LOB test flow (temporary object focus)

Large Inefficient XML Lookup

1 query performing 3 CROSS APPLY operations against 2 tables totaling 18GB resulting in 13 billion XML reader executions

Small inefficient XML lookup

1 query performing 3 CROSS APPLY operations against 2 tables totaling 600MB resulting in 475 million XML reader executions

Small inefficient XML lookups w/ delete - asynchronous

4 asynchronous queries using 3 CROSS APPLY operations against 8 tables @ 1.6GB resulting in 920 mil XML reader executions and final delete of specific records

Cursor operation in TempDB

1 query performing a series of data staging operations TempDB on 216MB EmailText data then utilizing a cursor to scan each record for a substring of text.

Long text string operation in TempDB

1 query loads 320MB of paragraph text into temp tbls then performs serialized scans for specific string values

GUID JOIN in TempDB

1 query loads 6GB of activity data across 3 temp tables, applies index on GUID column, then JOIN on column for final result set

CTE operation

1 query loads 2.5mil rows of messaging data into CTE using OUTER APPLY, then performs conversion select for result set

Large DDL & DML in TempDB

1 query loads 30mil rows @ 6GB of incident data into temp tbl then utilizes subqueries using a variety of grouping and sorting

LOB test flow (persistent object focus)

Large single remote extract

1 dataflow imports 30mil rows of incident data @ 6GB from remote server to physical table

Small single remote extract

1 dataflow imports 90k rows rows of xml data @ 1.3GB from remote server to physical table

Random READ

1k asynchronous queries selecting a random date range of records from 2 tbls totaling 600MB

Large DML operation

1 query loads 21mil rows @ 10GB of event detail data into physical table then performs distinct select on non-LOB values

Small multiple remote extract & complex agg

17 dataflows load 2GB of data from remote server into physical tables followed by 1 stored procedure performs several summary aggregations of 2GB incident topic data using physical tables, then returns grouped & sorted result set

Small inefficient XML lookups - synchronous

3 synchronous queries using 3 CROSS APPLY operations against 2 tables @ 1.3GB resulting in 890 mil XML reader executions

Maintenance Tasks

Diagnostic



DBCC CHECKTABLE on 26mil row 16GB table

Index Defragmentation



Rebuild 7 index partitions @ 3.3GB total

Database Restore



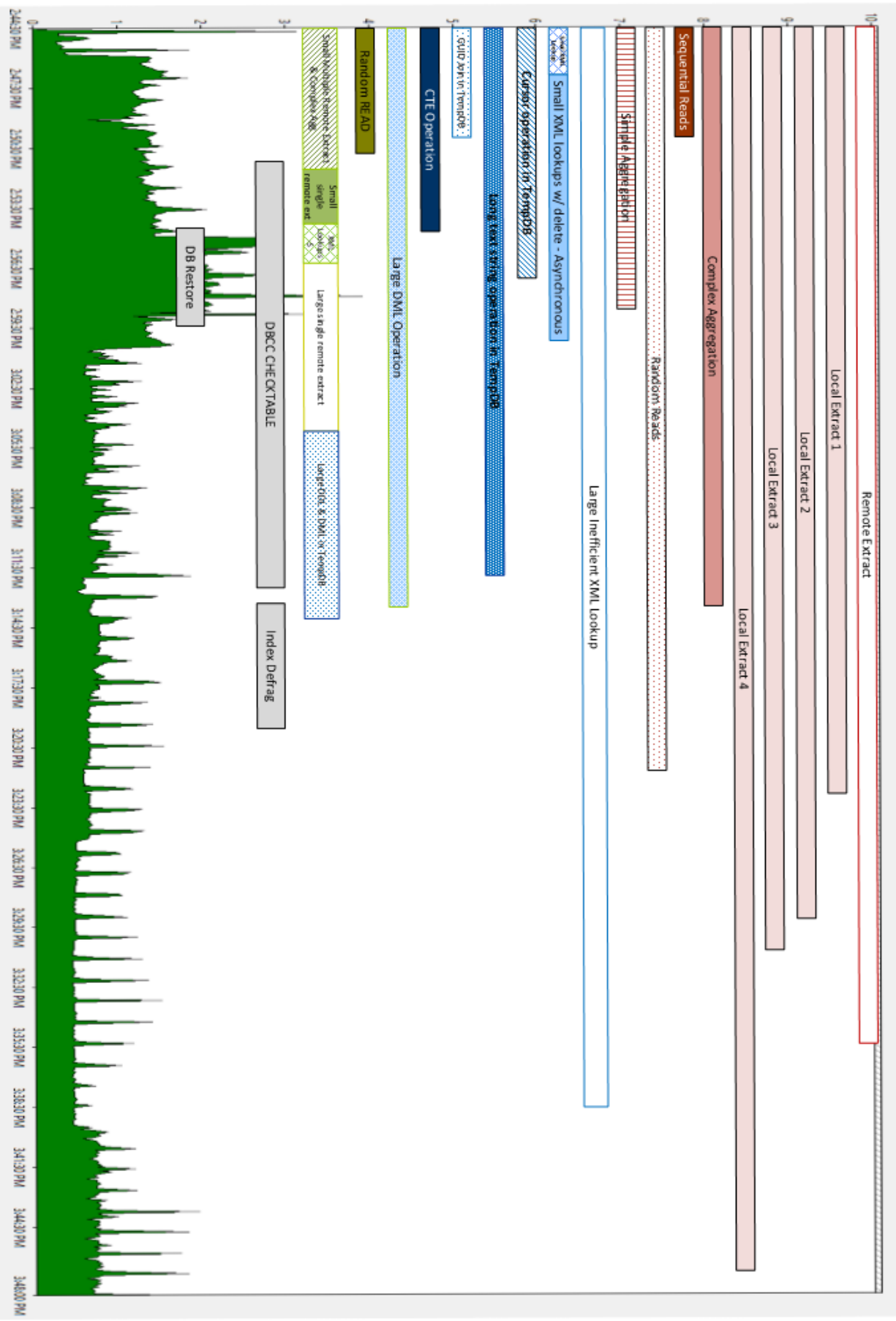
Restore 225GB DB

Success Criteria

Our success criteria was very simple - disk latency avg across all the XtremIO volumes should never exceeded 10 milliseconds and should be sustained below 5 milliseconds. Since it is common for most storage systems to experience short spikes in latency as cache is swapped, we allowed the transgression of the 5ms sustainment for activities shorter than 15 seconds.

Test Results

The avg test duration was 70 mins and during that time approx. 104k IOPS were generated from approx. 58GB of data movement. Disk latency across the XtremIO volumes never exceeded 4ms including the temporary spikes mentioned previously. The avg disk latency for the first phase of the test which was also the most I/O intensive was < 2ms. The avg disk latency reached its highest when the database restore executed causing sustained latency to bounce between 2ms & 3ms. Once the restore was completed, avg sustained latency remained < 1ms with isolated spikes at 1.75ms.



Summary & Conclusion

The success criteria outlined for the test was entirely met with final results actually exceeding our most optimistic expectations. Sustained latency rates were well below the recommended 10ms threshold for OLTP applications, therefore, to exceed this threshold in an I/O contentious BI environment is highly commendable.

As stated previously, the test package incorporated all the most I/O intensive operations that the CSSBI team will typically observe on the stand-alone DAP instance during peak usage periods. What's more, it incorporated maintenance operations that typically are schedule during non-peak usage periods. When considering how well the unit performed in an unlikely "worst case scenario" it is feasible that the I/O could be substantially increased further before any storage performance bottleneck would be encountered.

In addition, the fact that the XtremIO storage unit was able to easily handle the load with no care given to the actual placement of files on the volumes, or the segregation of files based upon READ/WRITE balance (ex. hosting TempDB files on the same volumes hosting LOB and clickstream data files) has the potential to greatly reduce the time and effort required of DBA's to organize the storage volumes and the files they host. SAN's achieve this to some degree with pools and block tiering, however, the setup of the pools and tiers still requires a SAN administrator. What's more, in most BI environments the rate of change and growth is exponential which requires DBA's to frequently work with the SAN administrator on storage optimization. However, the XtremIO storage product has the ability to be a "plug and play" solution where the optimization efforts of the DBA and SAN administrator are greatly reduced.

Despite the very large storage footprint of most BI systems, it's common that less than 20% of the data is actively used by the BI end-users. It's feasible that an XtremIO unit could be provisioned to host the most active data files, thereby allowing the storage of the remaining 80% on less expensive and highly redundant SATA pools either in DAS or SAN. Such a setup could ensure that storage performance was highly optimized without the higher associated expense of the hardware and the support personnel.

For more information:

<http://www.microsoft.com/sqlserver/>: SQL Server Web site

<http://technet.microsoft.com/en-us/sqlserver/>: SQL Server TechCenter

<http://msdn.microsoft.com/en-us/sqlserver/>: SQL Server DevCenter

<http://www.xtremio.com/> EMC XtremIO

Did this paper help you? Please give us your feedback. Tell us on a scale of 1 (poor) to 5 (excellent), how would you rate this paper and why have you given it this rating? For example:

- Are you rating it high due to having good examples, excellent screen shots, clear writing, or another reason?
- Are you rating it low due to poor examples, fuzzy screen shots, or unclear writing?

This feedback will help us improve the quality of white papers we release.

[Send feedback.](#)