

# TSM Paper – Replicating TSM

---

*(Primarily to enable faster time to recoverability using an alternative instance)*

Deon George, 23/02/2015

## Index

<b>INDEX</b>	<b>2</b>
<b>PREFACE</b>	<b>3</b>
<b>BACKGROUND</b>	<b>3</b>
<b>OBJECTIVE</b>	<b>4</b>
<b>AVAILABLE COPY DATA REPLICATION OPTIONS</b>	<b>4</b>
<b>TERMS USED IN THIS DOCUMENT</b>	<b>6</b>
<b>EXAMPLE SCENARIO USED IN THIS PAPER</b>	<b>7</b>
<b>FUNDAMENTAL REQUIRED RULES TO ACHIEVE ELECTRONIC COPY DATA REPLICATION</b>	<b>8</b>
<b>APPROACHES TO COPY DATA REPLICATION</b>	<b>9</b>
<b>TSM TO TSM REPLICATION WITH NODE REPLICATION</b>	<b>9</b>
<b>EXTERNAL METADATA (TSM DATABASE) REPLICATION</b>	<b>10</b>
TSM DATABASE REPLICATION VIA SAN	12
TSM DATABASE REPLICATION WITH HADR	14
<b>EXTERNAL DATA REPLICATION</b>	<b>15</b>
DATA REPLICATION VIA VTL	16
DATA REPLICATION VIA SAN (SYNCHRONOUS/ASYNCHRONOUS)	16

## Preface

Data Protection and Recovery is used for two primary purposes:

- For operational system recovery
- For business compliance data retrieval

Both of these purposes exist to primarily reduce financial impacts to a business as a result of not being able to perform these activities, but also exist to reduce a negative impact of an image of the business.

Since Data Protection and Recovery tools make secondary copies (aka “Copy Data”) of (then current) live business applications, it is recommended that a duplicate of this Copy Data be made to protect against the risk that the primary purposes cannot be achieved.

Some of the activities that affect the recovery of Copy Data for the primary purposes include:

- Media Failure
- Human Error
- Malicious Activities

Tivoli Storage Manager (TSM) is IBM’s Data Protection and Recovery product specialising in managing data for those two primary purposes (and more).

It’s architecture and core design, compliment a strategy to duplicate Copy Data to an alternative site, to provide recovery as a result of a (partial or full) failure in the primary site, including if the primary TSM server is impacted by the failure. This paper discusses those options.

This paper will refer to those options as Copy Data Replication.

It is also important, when choosing an external replication method, that the RPO (Recovery Point Objective) is well known and defined appropriately for the business. The value of the RPO will influence the final usage of the suggestions in this paper.

## Background

In the event of a catastrophic failure at the primary site and/or the primary TSM server, traditional disaster recovery practices would require restoring the TSM DATABASE, to a system with the TSM application installed (or pre-installed), making available the copy (or primary) media and activating the service for recovery.

The time for this approach is governed by the time it takes to:

- Provision a host to have the TSM application (if not available)
- Install the TSM application (if not pre-installed)
- Restore the TSM DATABASE (size vs time)
- Reconfigure TSM to recognise the available media, if the media now appears on alternative paths

Risk to this process includes:

- The age of the TSM DATABASE backup (often 24hrs old or older)
- The accessibility of the TSM DATABASE backup (if latest not available at the DR site, the DR DATA could be more than 1 day old)
- The ability to access the media (dependency on the source TSM server using technology that the target could use)
- The readability of the media, especially if the DATA on the media had not be “read/moved” recently (eg: long term archives), or if the media had been handled badly
- Correct configuration of TSM to support this approach (ie: configuration of re-use delay)
- The skills of the person executing the task

This paper will describe alternative approaches to improve this process, driven by requirements to reduce the risks and, still be a valid acceptable approach.

It has also been written to help design an implementation to cover un-scheduled outages. Since an un-scheduled outage often impacts business systems, the assumption here is that focus will be on returning an impacted business system back to operational status as quickly as possible.

Thus, in order to ensure that that task is successful, the infrastructure needs to be designed to reduce the affect of data loss, as a result of the unscheduled outage, and that any data loss is well known and defined.

## Objective

Provide:

- An electronic (thus little human intervention)
- Automated (thus occurs without human intervention)
- Continuous (so that RPO is less than 24hrs)
- Reliable (ie: usable when needed)
- Reduced time to first recovery (ie: , RTO is as low as possible)

## Available Copy Data Replication Options

Fundamentally, TSM (as are most data protection tools) is made up of two key components:

- The DATA in question
- The METADATA that represents the DATA in question

Typically, the METADATA describes which system owns it (or provided it to be stored), where it is currently stored, what it is called, and when restored, where it should be restored to (if not restored to an alternative location).

TSM provides several built-in Copy Data Replication techniques, with the technique relevant to this paper marketed as TSM Node Replication. This technique already replicates the DATA and METADATA using a TSM to TSM replication technique.

When TSM Node Replication is not used, external methods that result in the DATA being replicated and the METADATA being replicated can be used.

The available methods discussed in this paper to replicate the METADATA are:

1. DB2 HADR
2. SAN Device replication

If one of the above methods is used, then it should be used in combination with one of the following methods to replicate the DATA.

1. VTL replication
2. SAN Device replication

Not considered in this paper

- METADATA replication and DATA replication, using file level copy techniques – for example RSYNC. While technically this capability could be used, functionality it is inferior to one of the methods above.

IBM would recommend that TSM Node Replication be used in lieu of this approach.

- METADATA backup to a device, where that device can replicate to the target site. For example, a TSM Database Backup to the VTL, where the VTL replicates to a remote site, or a METADATA backup to a file system, where the filesystem LUNs are replicated by a the storage controller.

While this is perfectly acceptable method to get the METADATA and DATA to a remote site, it is not discussed in this paper.

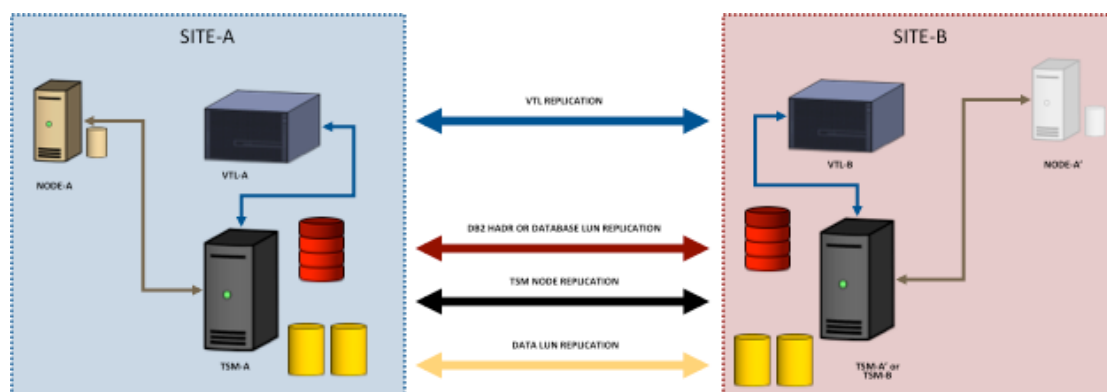
- Data replication where the DATA is stored in TSM (random access) “DISK” pools. This paper discusses DATA that is stored in sequential access pools, for example “FILE” or “VTL” device classes.

Sequential access pools would be preferred over random access disk pools, providing several points of control to ensure that the DATA and METADATA are available with the appropriate in age at the remote site.

## Terms used in this document

- DATA – The files, application data that TSM manages on behalf of Data Protection Clients.
- DATABASE – The TSM database, comprised of METADATA, that references DATA stored in a STORAGE POOL
- DISK – DATA is stored on a DEVICE TYPE of DISK
- DEVICE TYPE – The type of device used to store DATA (from TSM's view of the world).
- FILE – DATA is stored on a DEVICE TYPE of FILE, which are physically files on a file system accessed by the TSM server application.
- METADATA – The entries in the TSM database that describes the DATA stored in a STORAGE POOL, therefore making it possible to retrieve the DATA when requested.
- STORAGE POOLS – The pool of storage that DATA is stored on by the TSM server (from TSM's view of the world).
- VTL – DATA is stored in an external device, placed there either by the TSM server, or by a TSM storage agent. The METADATA is managed by TSM.

## Example scenario used in this paper



Data Protection client NODE-A, primarily sends its DATA to TSM-A, located at SITE-A. The TSM server is configured to store this backup either externally in a VTL VTL-A (and thus the VTL may further process the DATA), or stored internally by TSM in a storage pool that uses the FILE device class.

At the Disaster Recovery site SITE-B, there are separate physical infrastructure components that will be used to instantiate a TSM server and enable the recovery of DATA from that storage infrastructure.

### In the case of TSM Node Replication:

The TSM server would be called TSM-B and it will manage copies of DATA from TSM-A on behalf of TSM-A.

Actually, TSM-B will treat NODE-A as its own, store NODE-A's DATA as per TSM-B's configuration, and should NODE-A ever directly connect to TSM-B, will hand back NODE-A's DATA on request.

(TSM-B would only get DATA owned by NODE-A from TSM-A, not from the node directly.)

TSM-B may also have its own workload to manage, that TSM-A is not aware of.

### In the case of External Replication:

The TSM server would be a clone of TSM-A at SITE-A, and thus is called TSM-A'. This server will behave as if it was the original TSM-A server at SITE-A, and thus would not manage any other data. (In fact, it is not aware that it is a clone.)

Since it is a clone of TSM-A, it is essentially not functioning until specifically promoted to be TSM-A, and at which time the original TSM-A may be (and possibly should be) offline.

The external storage device (VTL or replicated LUNs) has identical copies of the data at SITE-A, and depending on the functionality of the storage device, may also hold other data. This other data is unknown (and irrelevant) to TSM-A'.

## Fundamental Required Rules to achieve Electronic Copy Data Replication

### In the case of External Replication:

If, at the remote site, replication is occurring externally to TSM, the synchronization of DATA and METADATA is critical, and where it cannot be guaranteed to arrive at the target at the same time,

1. The DATA must be at the target before the METADATA
2. The DEVICE TYPE used is sequential volume (ie: FILE/VTL)
3. The "REUSEDELAYDELAY" parameter enabled on the STORAGE POOL. (This parameter, in days, must be greater than the maximum age difference between the DATA and METADATA at the target.)

Implied here, is that access to Data Protection data may not provide the latest copy of DATA, if there is a difference between the age of the DATA and METADATA.

If there is a known scenario, were the DATA replication may be stalled, and the METADATA replication could continue (for example 1 of 2 communication links could be affected by an outage), then it is also a requirement, that additional techniques be implemented to provide an ability to force the METADATA to a prior time point so when it is used, it is older than the (current) DATA at the target site.

One technique that can be used to achieve this is FlashCopy (aka Snapshot) of the DATABASE, which could be done in the Storage Infrastructure itself (LUN Flash Copy), or by the host at a Software layer, for example, LVM snapshot.

This technique should run automatically and regularly, so that a reasonable recent RPO can be achieved given the age of the DATA at the time.

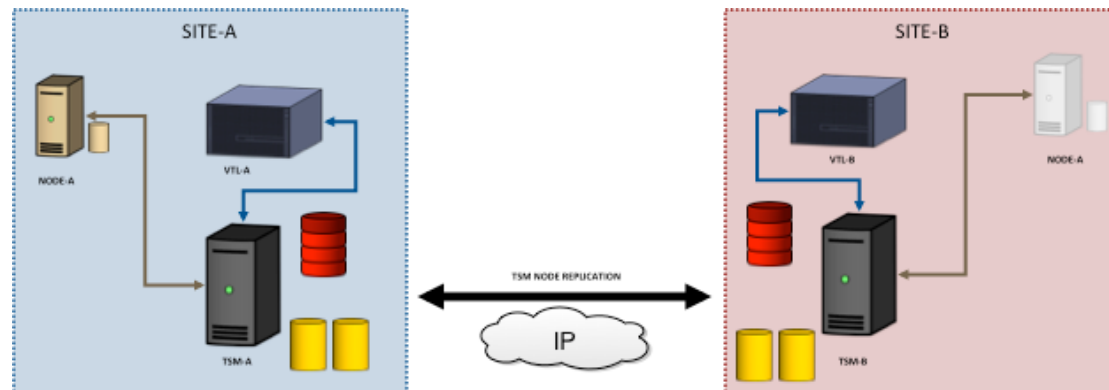
### In the case of TSM Node Replication:

Where replication is occurring internally with TSM (using TSM's built-in feature Node Replication), then the synchronisation requirement is satisfied, as it is managed internally by TSM itself.



## Approaches to Copy Data Replication

### TSM to TSM Replication with NODE REPLICATION



This is the most efficient and most reliable method to use. It is also the recommended approach.

Why?

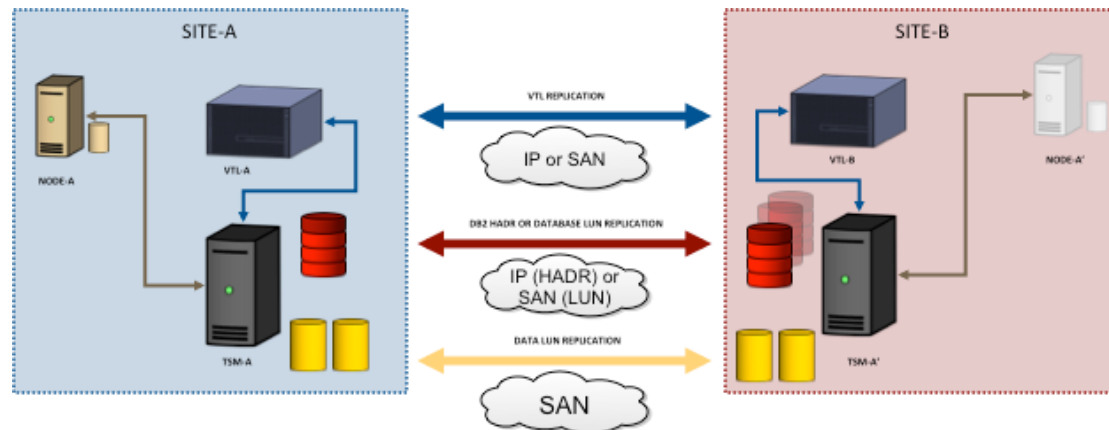
- Synchronisation is managed by TSM, and thus therefore no longer critical
- Most efficient in terms of network bandwidth required:
  - Especially when de-duplication is leveraged, together with client compressed data.
  - Replicates only the required items from source to target
- No ongoing configuration required to use the data at the target
- The target is “Hot”, IE: The target is open and ready to service requests, and thus can restore data at any time, as often as required, including when the source site is still operational
- With TSM 7.1.1, the target can contain more or less versions of data than the source and the source and target can apply different retention times
- With TSM 7.1.1, the target can be used to recover the source, in the event that the source has data loss.
- It is TSM application to TSM application, which implies integrity both at the source and at the target
- Source TSM platform and Target TSM platform can be different Operating Systems and hardware configurations
- Source and target TSM servers can use storage pool hierarchies comprising different storage devices (whereas device replication is restricted to a single storage device and the same device must be used at source and target)
- Source TSM version and Target TSM version can be N+1 (**double check?**)
- Clients automatically failover to use the target for recoveries (Currently the BA Client Only)

Notes:

- Only available via the IP network

Where IP to IP replication is available, and DATA is not stored in VTL, this is the preferred and recommendation approach to use.

### External METADATA (TSM Database) replication



These approaches are available where an IP to IP approach cannot be used (for both DATA and METADATA), or it is preferred to use a storage device level replication capability in lieu of TSM's NODE REPLICATION.

Additional Recommendations for these approaches (unless the DATA and DATABASE (including logs) are in the same consistency group):

- Since the DATA replication is performed autonomously via another process, and potentially a different communications link, that the target's DATABASE is frequently "snap shot" to obtain different consistent points in time.

This requirement exists, to ensure that the METADATA can be used and be older than the corresponding DATA that has been replicated by a different process.

If a there is no mechanism to perform a regular (and automatic) "snap shot" one of two actions must occur before the TSM instance can service recoveries:

- An AUDIT VOLUME
- Delay using the METADATA until it is known that the DATA synchronisation has completed, and the corresponding DATA is now newer than the METADATA that references it.

For example,

- TSM-A's database (the METADATA) is being replicated by an (a)synchronous process from SITE-A to SITE-B,
- Via a different communications link, the DATA is being replicated by a VTL from SITE-A to SITE-B.
- At 10:05 am, the DATA replication link fails, however, the METADATA replication link is not affected by the failure.

- At 11:00 am, SITE-A fails, and it is declared to promote the DR site to production status, which requires activating TSM-A'

At 11:00 am, at SITE-B, the TSM-A's DATABASE, could include Data Protection activity from 10:05 am to 11:00 am, however, the VTL only has data up until 10:05 am.

If the TSM server is activated with the 11:00 am DATABASE, it may contain references to data, that currently is not available on the VTL at SITE-B, or it may contain references to data that was moved, but the VTL at SITE-B has no knowledge of that movement.

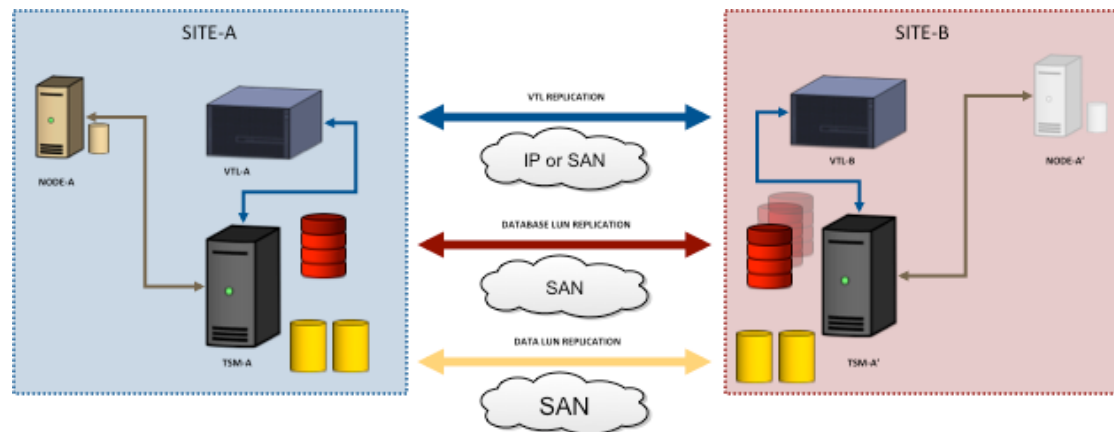
If a Flash Copy/Snapshot technique was used for the DATABASE at SITE-B, and was taking snapshots every hour (on the hour), it could be possible to revert the SITE-B's version of the TSM DATABASE back to 10:00 am and thus the TSM DATABASE's would correctly be able to access all data in the VTL referenced by the 10:00 am DATABASE snapshot.

Any data that the VTL received between 10:00 am and 10:05 am would be unknown to that version of the TSM DATABASE, and would not be recoverable.

Depending on the severity and impact of the unscheduled outage that promoted the disaster recovery site to product status, the inability to access this 5 minutes of data may be negligible.

(Should it not be negligible, more frequent snapshot points could be used to mitigate this concern.)

## TSM Database Replication via SAN



This configuration is where the TSM Database is stored on a SAN Storage Device, and that device provides the ability to replicate the LUNs to a remote SAN Storage Device.

TSM makes no recommendations nor a requirement on what SAN Storage Device is used. The requirement on SAN Storage Device is the responsibility of the Operating System (OS), in that the OS must support the attachment of the vendor SAN Storage Device, and may require specific drivers to do so.

Some important notes to be aware of with this configuration:

- In terms of preference, performance and operational optimisation, it is always recommended to use SSD (or similar technology) for the TSM DATABASE, which may impact the use of this replication capability, unless the use SSD (or similar technology) is available via the SAN Storage Device.
- Where it is possible to have DATA volumes and METADATA volumes (including database log volumes) in a “SAN device consistency group”, that approach should be adopted to mitigate any synchronisation requirements.

Implied with this approach is:

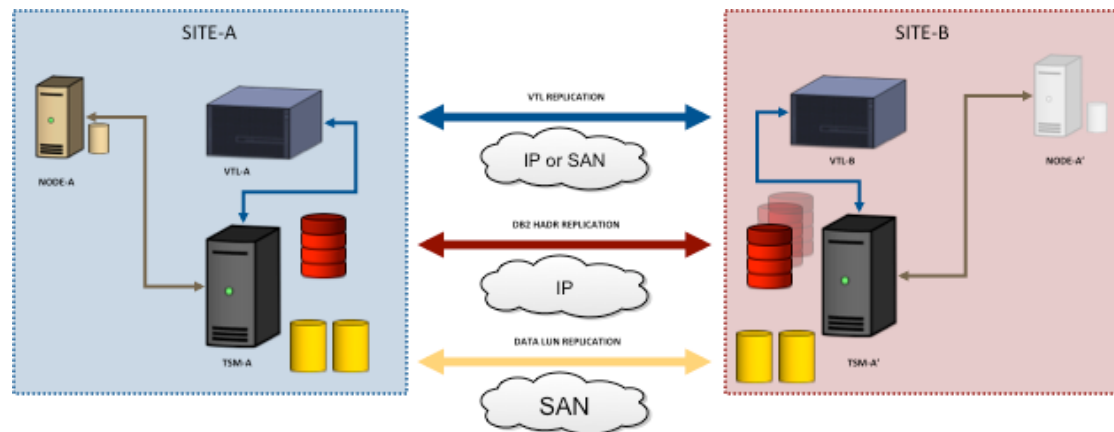
- Replication is external to TSM, and thus synchronisation of DATA and METADATA is critical if a “SAN device consistency group” is not used.
- Integrity of the TSM DATABASE is critical, and thus TSM Database and Database Logs need to be included in the replication, to ensure a crash consistent recovery is possible.
- The target server’s view of the world is that the TSM has not been shut down cleanly when it is started, and thus it will be in a state requiring crash consistent recovery.

This approach is managed by the SAN Storage Device infrastructure, and the ability to monitor its status, understand the time delay between the primary and disaster recovery sites and rectify any problems is with the SAN Storage Administrator.

The act of promoting the disaster recovery site's instance to production will require:

- Suspension of the replication (if not done already), and
- Enable the target volumes to be used for read-write operations

## TSM Database Replication with HADR



This configuration is where the TSM Database (DB2) is replicated using a DB2 capability known as DB2 HADR. The basic architecture of HADR is simple. It is based on the shipping of log records from the primary database to a standby database.

DB2 HADR provides several replication options, and any option could be used to achieve the METADATA replication requirement. (On IBM's wiki, it references using SYNC mode. This wiki was authored when TSM was using an older version of DB2 and the other modes were not available.)

The options are:

- SYNC – In this mode, writes are not considered successful until successfully applied to both the primary and target instances and recorded in the log files of both instances.

This mode is not typically used by TSM, as it impacts the performance of TSM.

- NEARSYNC – Similar to SYNC, instead the write is considered successfully when only applied to memory in the target instance.

This mode is not typically used by TSM, as it impacts the performance of TSM.

- ASYNC – In this mode, writes are only considered successful when applied to the primary instance, and have been delivered to the network (for the remote site).

This mode is not typically used by TSM, as it impacts the performance of TSM.

- SUPERASYNC – In this mode, writes are not impacted by the performance of the link or the remote site, nor the availability of the remote site.

This would be the preferred mode to be used with TSM, as the performance of the link to the remote site, as well as the performance of the remote system does not impact the primary TSM server.

To assist with forcing the DB2 database to be behind the DATA replication, DB2 provides HADR Replay Delay parameter “`hadr_replay_delay`”, which delays the applying of changes received from the primary database.

This parameter, in seconds, could be used to ensure that METADATA is always behind the DATA. It should be used with the “`hadr_spool_limit`” parameter, please refer to the DB2 documentation on the usage of these parameters.

Implied with this approach is:

- Replication is external to TSM, and thus synchronisation of DATA and META DATA is critical.

The setup and configuration of HADR is outside the scope of this document. There are Redbooks and/or the DB2 documentation that can be used to assist with the setup and configuration of DB2 HADR.

### External DATA replication

Both these scenarios assume that TSM stores the DATA on the storage device using a sequential media format (ie: device type of FILE, or a tape device type when using a VTL).

It is recommended in the TSM Storage Pool configuration, that the REUSEDELAY parameter be used, and at a minimum be set to 1 (ONE) to ensure that volumes used are not overwritten for at least one day after all files have been deleted from the volume. This provides protection should the replication time point of the METADATA be later than that of the DATA.

The REUSEDELAY parameter should be set to the maximum expected time difference between replication of the METADATA and DATA to the target.

As a side effect of utilising this parameter, storage that would be freed, as a result of TSM moving data out of those volumes, will not be actually freed until expired by the delay imposed by the REUSEDELAY parameter.

The REUSEDELAY parameter is important for many scenarios, here is an example of 1 scenario, where the REUSEDELAY parameter protects data from the remote site from being deleted too soon:

- TSM-A’s database (the METADATA) is being replicated by an (a)synchronous process from SITE-A to SITE-B,
- Via a different communications link, the DATA is being replicated by a VTL from SITE-A to SITE-B.

- At 10:05 am, the METADATA replication link fails, however, the DATA replication link is not affected by the failure.
- At 10:30 am, while the primary TSM server TSM-A is still active, it completes a RECLAMATION process that removes all data off of VOLUME “VOL1”, and stores that data on “VOL99”.
- The data movement that has occurred between “VOL1” and “VOL99” is successfully replicated to the remote site.
- At 11:00 am, SITE-A fails, and it is declared to promote the DR site to production status, which requires activating TSM-A’

When TSM-A’ is promoted to production status, the DATABASE at SITE-B is only aware of transactions that occurred up to 10:05 am. It has no knowledge of the data movement that moved remaining data off of VOL1 to VOL99.

However, since the VTL replication was unaffected by the outage, the VTL at SITE-B is aware of the data movements.

If the REUSEDELAY parameter was not enabled, the act of emptying VOL1 may result in it being overwritten by the primary VTL, and thus being overwritten in the remote VTL. Any attempt to restore files that were on that volume (as at 10:05 am version of the DATABASE), would result in errors, since that volume no longer exists or has been overwritten.

Setting the REUSEDELAY parameter to 1 (one) day, will ensure that the VOL1 is not actually overwritten until 1 day has passed, the DATA on VOL1 at SITE-B would still be recoverable.

(When TSM moves data, it doesn’t remove it from the source, and write it to the target, it de-references from the source and writes a copy to the target. Thus restoring the DATABASE to “the past”, will result in those de-references not existing (yet), and the original references still in place.)

### Data Replication via VTL

Tivoli Storage Manager makes no requirements on how VTL replication is performed, other than to ensure that the VTL replication is always ahead of the METADATA replication.

Where it is possible that VTL replication can be behind METADATA replication, then additional Flash Copy/Snapshot techniques should be used on the METADATA, so that it can be forced to an older version, than that of the replicated DATA.

NOTE: Some VTLs may not be able to replicate data until the mounted volume has been dismounted (or closed). This can be achieved in TSM, by setting the MOUNTRETENTION value to 1 min (it’s default value is 60 minutes).

### Data Replication via SAN (synchronous/asynchronous)



Tivoli Storage Manager makes no requirements on how SAN replication is performed, other than to ensure that the SAN replication of the DATA is always ahead of the replication of the METADATA, if a consistency group cannot be used.

Where it is possible that VTL replication can be behind METADATA replication, then additional Flash Copy/Snapshot techniques should be used on the METADATA, so that it can be forced to an older version, than that of the replicated DATA.

Where it is possible to include the DATABASE LUNs in the same consistency group as the DATA LUNs, then this should be used, and mitigates the requirement to ensure that the METADATA is older than the DATA.