

zHyperLink Write Support

zHyperLink technology now provides update write capability for both simplex and metro mirror environments.



By: Tariq Hanif, Brian Lee, Dale Riedy Published: November 19th, 2021 Read time: 5 minutes

Overview

I/O latency as viewed by the application is more than just I/O service time. It also includes interrupt processing time, time waiting to get dispatched on a CPU after an I/O has completed, and delays caused by having to reload the CPU cache. Traditional I/O to DASD is asynchronous, which means that the application requests a record, the operating system initiates an I/O request, and the application is suspended until the I/O completes and the application is redispatched on a CPU. zHyperLink I/O is synchronous, which means that the CPU spins waiting for the I/O request to complete. This eliminates the delays normally associated with asynchronous I/O.

I/Os that are good candidates for zHyperLink are those where the application is suspended until the I/O completes. Two examples are Db2 synchronous database reads that occur when the application requests a record that is not in the buffer pool and Db2 log writes issued at transaction commit time. I/Os that are not good candidates for zHyperLink are reads used for prefetching or read-ahead and delayed writes because both are performed asynchronous to the execution of the application.

zHyperLink write requests have an additional requirement: the write I/Os must follow a log write pattern, which means that the writes must always go forward in the data set without skipping any records. The same record may be written multiple times, but the writes never go backwards. When you are finished with the last record in the data set, you can restart at the beginning of the data set.

In 2019, zHyperLink write support was introduced for simplex and local Metro Mirror environments (see Figure 1). In Metro Mirror configurations all DS8000s must be within 150m of the Z server and connected via zHyperlink.

In 2020 for z/OS V2R3 and above, zHyperLink write support was provided for DS8000 Global Mirror (asynchronous mirroring) configurations (see Figure 2 for an example). This does not include support for Extended Remote Copy (XRC). Db2 for z/OS currently exploits zHyperLink writes for active log data sets. Media Manager support for dual logging provides the ability for Db2 to issue zHyperLink writes to two active log data sets simultaneously¹.



Figure 1: zHyperLink writes in a Metro Mirror environment.



Figure 2: zHyperLink writes in a Global Mirror environment.

Hardware and I/O Configuration Requirements

zHyperLink is a point-to-point connection between the IBM Z processor and a DS8000. The maximum distance is 150 meters. To use zHyperLink, you need one or more zHyperLink Express features (cards) in the IBM Z processor. Each feature supports up to 2 ports and up to 16 features (32 ports) can be installed.

On the DS8000 side, each zHyperLink connection plugs directly into a port in the I/O enclosure; a special host adapter card is not required. Up to 2 connections per I/O enclosure is supported. The maximum number of connections supported is model dependent and is based on the number of enclosures and cores. zHyperLink connectivity is also required to the secondary for write support in Metro Mirror

environments, which means both primary and secondary must be within 150m of the IBM Z processor.

FICON connectivity to the DS8000 is still required for initialization for performing I/Os that are not eligible for zHyperLink and as a fallback if the zHyperLink I/O fails. Using zHyperLink does not require you to reduce the maximum number of FICON channel paths configured for a device.

In the IBM Z I/O configuration, define the PCIe function definitions (PFIDs) for the zHyperLink ports. The association between the zHyperLink port and DS8000 is discovered dynamically. Define 4 PFIDs per zHyperLink port for each logical partition that is sharing the DS8000.

Software Requirements

Find the PTFs for zHyperLink by searching for Fix Category IBM.Function.zHyperLink, APAR keyword HYPERL/K.

Enabling zHyperLink Write on z/OS

On a z/OS LPAR, zHyperLink is disabled by default. To enable zHyperLink write processing during IPL time, include the following in the IECIOSxx parmlib member²:

ZHPF=YES, if zHPF is not already enabled on the LPAR.

ZHYPERLINK,OPER=WRITE or OPER=ALL. If ALL is defined, both read and write processing is enabled.

zHyperLink write can also be enabled dynamically on a z/OS LPAR using these z/OS commands³:

SETIOS ZHPF=YES, if zHPF is not already enabled on the LPAR.

SETIOS ZHYPERLINK, OPER=WRITE or OPER=ALL.

To verify if zHyperLink support is enabled, issue:

- DISPLAY IOS, ZHYPERLINK to display whether zHyperLink is enabled for writes for the system (see Figure 3).
- DISPLAY M=DEV(devno) to display whether zHyperLink is enabled for writes for a device (see Figure 4).
- DISPLAY M=DEV(devno),ZHYPERLINK to display whether zHyperLink is enabled for writes for a device. If zHyperLink is disabled, the reasons why zHyperLink is disabled are displayed.
- DISPLAY M=CU(cuno) to show whether zHyperLink is enabled for writes at the control unit level and to show the state of the PFIDs (see Figure 5).

D IOS,ZHYPERLINK IOS634I 16.12.16 IOS SYSTEM OPTION Ø61 Zhyperlink IS enabled for read and write operations

Figure 3: DISPLAY IOS, ZHYPERLINK command output.

```
M-DEV(@A7@F)
IEE174I 16.18.33 DISPLAY
                          M 296
DEVICE ØA7ØF
               STATUS-ONLINE
                       50
CHP
                            52
                                  69
                                       6В
ENTRY LINK ADDRESS
                       35
                             1235
                                  02
                                       02
    LINK ADDRESS
                             1207
                       07
                                       08
                                  08
   H ONLINE
   PHYSICALLY ONLINE
  TH OPERATIONAL
                             Y
MANAGED
                            Ν
                                  Ν
                                       Ν
  NUMBER
                       A701
                            A701
                                  A701
                                       A701
INTERFACE ID
                            0033
                                  0031
                                       0231
                       0230
MAXIMUM MANAGED CHPID(S)
                          ALLOWED:
                                     0
DESTINATION CU LOGICAL ADDRESS 🗕 67
                   = 002107.996.IBM.75.0000000LBL41.0230
SCP
   CU ND
    TOKEN NED
                     002107.900.IBM.75.0000000LBL41.6700
SCP
   DEVICE NED
                     002107.900.IBM.75.0000000LBL41.670F
                     500507630AFFD03A
WWNN
 PERPAV ALIASES CONFIGURED -
                               16
ZHYPERLINKS AVAILABLE = 4
FUNCTIONS ENABLED - MIDAW,
                            ZHPE, HS, XPAV,
                                              ZHYPERLINK
```

Figure 4: DISPLAY M=DEV(devno) command output.

```
) M=CU(X701)
EEE174I 16.10.19
CONTROL UNIT A701
                        DISPLAY
                                        055
     LINK ADDRESS
Physically on
                     ONLINE
       VALIDATED
     GED
       - CHPID
     - CU INTERFACE
RFACE ID
                                               0031
           MANAGED CHPID(S)
                                    ALLOWED
                                                   0
                 CU LOGICAL
          TION
                                              IBM.75.0000000LBL41.0230
                             002107.
                          - 002107.996.IBM.75.0000000LBL41.
- 002107.900.IBM.75.0000000LBL41.
 OKEN NED
                                                                           6700
                                     7630AFFD03A
                             ZHPF, 0500,
0300, 0500,
1F01, 2301,
5600, 7300,
5901,
   NCTIONS ENABLED
      CU PEERS
                                                                                            1001
                                               2501,
7700,
                                                        2801,
7900,
                                                                 2A01,
7800,
                                                                          3101
                                                                                   3200
DEFINED DEVICES
   0A700-0A70F
FINED PAV ALIASES
      LE HYPERPAV ALIASES - 16
                           Port
                                    Linkid
                 0160
                                     0080
                                                 Alloc
                                                               Oper
                 0160
                                     0080
                                                 Alloc
                            Øl
                                     0080
                            Ø1
                            øı
                                     0180
                  0100
                            øı
                                         80
```

Figure 5: DISPLAY M=CU(cuno) command output.

Enabling zHyperLink Write on Db2

- zHyperLink write support is available on Db2 12 for z/OS.

To enable zHyperLink, use the Interactive Storage Management Facility (ISMF) interface to change the zHLink Write parameter of the storage class used to allocate the Db2 active log data sets to YES from the default of NO. Use the VARY SMS command to override the storage class setting for SMS managed data sets or to allow zHyperLink writes to be used for non-SMS managed data sets. Note: The VARY SMS command change does not persist across an IPL. To exploit zHyperLink support for Db2 active log datasets, the ZHYPERLINK Db2 subsystem parameter must be set to either ACTIVELOG or ENABLE. This can be dynamically updated while Db2 is running by issuing the Db2 SET SYSPARM command⁴

Enabling zHyperLink Write on DS8000

- zHyperLink write support for synchronous mirroring is available on the DS8880 and DS8900.
- zHyperLink write support for asynchronous mirroring is available only on the DS8900.

To enable zHyperLink, logon to DS8000 HMC as an administrator, select Settings and then System to see the list of supported licensed functions. The zHyperLink Function shows the options to enable I/O Write and I/O Read (see Figure 6). Enabling zHyperlink writes on the DS8000 will result in a quiesce/resume of the DS8000 servers similar to a code load, and so while this is non-disruptive, it would typically be scheduled for a quiet time.

IBM	IBM DS8000	IBM.2107-75LBL41 System
ᢙ	Dashboard	Licensed Functions zHyperLink
¶ م	Monitoring	Easy Tier zHyperLink point-to-point connections provide ultra-low latency for selected eligible I/Os.
₿	Pools	zHyperLink I/O Write V Enabled I/O Read V Enabled
8	Volumes	Date and Time
	Hosts	Advanced
H	Copy Services	
<u> </u>	Access	
ŵ	Settings	

Figure 6: Enabling zHyperLink on DS8000.

For Metro Mirror environments, the devices on which the Db2 log data sets reside must be in the duplex state, in a HyperSwap configuration, and zHyperWrite must be enabled. In the IBM test lab, the solution was tested when the devices were in a simplex and the following mirroring configurations:

- Single site
 - Simplex
- 2-site
 - Metro Mirror (MM) (see Figure 7)
 - Global Mirror (GM) (see Figure 8)
- 3-site
 - Multi-Target (MT) MM/MM
 - Multi-Target MM/GM (see Figure 9)
 - MGM (see Figure 9)



Figure 7: Two Site Metro Mirror environment.



Figure 8: Two Site Global Mirror environment.



Figure 9: Three Site Multi-Target MM/GM or MGM environment.

zHyperLink Write Performance Results

zHyperLink writes provide lower latency than zHPF in all replication configurations. Figure 10 compares the performance of zHyperLink writes to zHPF writes with no mirroring and with different mirroring environments (Metro Mirror, Global Copy, Global Mirror). These results were achieved with 26-meter zHyperLink connections using a workload that simulates Db2 log writes that were 4K in size. Note: In this test, zHyperWrite was not enabled when using zHPF in a Metro Mirror environment so the response time shown is higher than it would have been with zHyperWrite enabled. With zHyperWrite, the response time for Metro Mirror would be similar to no mirroring.





Figure 11 compares zHPF to zHyperLink performance in a Global Mirror environment with a mixed random/sequential read and write workload and using a program that emulates Db2 log writes. zHyperLink reads were not enabled for this test. There was a 4x reduction in log write response time with zHyperLink. Overall, response time showed only a small improvement because reads made up a majority of the I/Os. In addition, this test showed that the success rate for zHyperLink writes was unaffected by Global Mirror and the consistency group formation time, which affects the Recovery Point Objective (RPO), was not affected by zHyperLink writes (1.1 seconds with and without zHyperLink).



Figure 11: zHyperLink write vs FICON with Global Mirror consistency group.

zHyperLink Write Statistics and Reports

RMF shows synchronous I/O performance data at the DASD device and PCIE function level. Figure 12 shows an example of the RMF Synchronous I/O Device Activity report, which contains the zHyperLink read and write activity as well as the asynchronous I/O activity.

			S	УNС	HRONO	US I,	O DE	VIC	E A C	тіVІ	тү				
z	/OS V2R4			SYSTEM RPT VEH	ID SY1B RSION V2R	4 RMF	STAR END	T DATE DATE	10/30/2 10/30/2	020-14.1 020-14.2	.9.00 25.00	INTERVA CYCLE 1	L 000	. 05 . 59 SECONDS	3
TOTAL SAMPLES	= 360	IODF =	1B	CR-DATE	E: 07/16/	2020 CR-	TIME: 0	7.22.2	1 AC	T: ACTIV	ATE	s.	9	\$	8
STORAGE DEV	DEVICE	VOLUME	LCU -	SYNCH	4 I/O	ASYNCH	-SYNCH	1/0 -	ASYNCH	TRANSFE	R RATE	REQ	LINK	CACHE	REJECTS
GROUP NUM	TYPE	SERIAL		READ	WRITE	1/0	READ	WRITE	I/O	READ	WRITE	SUCCESS	BUSY	MISS	READ WRITE
DC1NLOG 0B024	3390A	TDB024	0050	0.000	5482.85	0.547	0.000	0.064	0.786	0.000	109.9	100.0	0.00	0.00	0.00 0.00
		LCU	0050	0.000	5482.85	0.547	0.000	0.064	0.786	0.000	109.9	100.0	0.00	0.00	0.00 0.00
DC1NLOG 0B124	3390A	TDB124	0051	0.000	5482.85	0.558	0.000	0.072	0.815	0.000	109.9	100.0	0.00	0.00	0.00 0.00
		LCU	0051	0.000	5482.85	0.558	0.000	0.072	0.815	0.000	109.9	100.0	0.00	0.00	0.00 0.00

Figure 12: RMF synchronous I/O device activity report.

The DISPLAY LOG command for Db2 for z/OS shows if zHyperLink writes is enabled for

the active logs and whether the last log write for the active log data sets used zHyperLink (see Figure 13).

-DAP0 DIS LOG										
DSNJ370I -DAP0 DSNJC00A LOG DISPLAY 188										
CURRENT COPY1 LOG - DBA1.DAP0.LOGCOPY1.DS03 is 82% FULL										
CURRENT COPY2 LOG - DBA1.DAP0.LOGCOPY2.DS03 is 82% FULL										
H/W RBA - 00000002C14FE356EE8										
H/O RBA – 00000002C149640AFFF										
FULL LOGS TO OFFLOAD - 0 OF 6										
OFFLOAD TASK IS (AVAILABLE)										
SOFTWARE ACCELERATION IS ENABLED										
ZHYPERLINK WRITE IS ENABLED										
LAST LOG WRITE FOR COPY1 USED ZHYPERLINK										
LAST LOG WRITE FOR COPY2 USED ZHYPERLINK										

Figure 13: Db2 DISPLAY LOG command output.

Alternatively, DFSMS collects detailed zHyperLink write statistics and makes them available with the SMF type 42 subtype 6 records and the DISPLAY SMS,DSNAME(dsn),STATS(ZHLWRITE) command (see Figure 14). The statistics are provided at both a Db2 request level as well as the device level.

D SMS, DSNAME (DBA1.DAP0.LOGCOPY1.DS02), STATS (ZHLWRITE) IGW289I 419 D SMS, DSNAME, STATS (ZHLWRITE) Start of Report										
DATA SEI DBAL.DAPU.LOGCO.	$\begin{array}{c} \text{PII.DSUZ.DATA} \\ \text{OO1} 12.51.20 1 \end{array}$	0166								
STATISTICS Since 05/20/2	021 13:51:39.14	12100								
SUMMARY										
TOTAL %SYNC%ASYNC										
WRITE REQUESTS WRITES	SKIP LNKBSY	^EST	MISC	DISABL						
751121 99.89	0.03 <0.01	<0.01	<0.01	0.00						
%ASYNC										
	MISS DELAY	DUAL								
	<0.01 0.04	0.01								
DEVICE STATISTICS										
TOTAL %SYNC%ASYNC%										
SSID DEVNO WRITES WRITES	SKIP LNKBSY	^EST	MISC	MISS	DELAY					
1241 0030F 106e04 99.96	0.01 <0.01	<0.01	<0.01	<0.01	0.01					
12E1 08C64 106e04 99.95	0.01 <0.01	<0.01	<0.01	<0.01	0.01					
1241 00305 106e04 99.94	0.01 <0.01	<0.01	<0.01	<0.01	0.01					
D SMS, DSNAME, STATS (ZHLWRITE) End of Report										

Figure 14: DISPLAY SMS command output.

Omegamon XE for Db2 Performance Expert shows the average elapsed time for class 3 suspensions for synchronous I/O log writes⁵. The TIME/EVENT for SYNCHRONOUS LOG WRITE I/O shows the benefit of using zHyperLink writes for Db2 active logs when compared to reports prior to enabling zHyperLink writes (see Figure 15).

LOCATION: GROUP: MEMBER: SUBSYSTEM: DB2 VERSION:	DSNC1N0 DSNC1N0 DC1N DC1N V12		OMEGAMON XE FOR DB2 PEF ACCOUNTING F ORDER: PRIN SCOPE:	RFORMANCE EXPER REPORT - LONG MAUTH-PLANNAME MEMBER	RI	PAGE: 1-1 EQUESTED FROM: 10/30/ TO: 10/30/ INTERVAL FROM: 10/30/ TO: 10/30/	20 14:19:00.00 20 14:25:01.01 20 14:19:00.00 20 14:25:01.00	
PRIMAUTH: USR ELAPSED TIME	T031 PLANNA DISTRIBUTION	ME: db2jcc_a		CLASS 2	TIME DIST	RIBUTION		
APPL ===== DB2 ===== SUSP =====	> 10% >	> 32%	> 58%	CPU SECPU NOTACC SUSP	====> ===> 6% ======	17%	> 77	8
AVERAGE	APPL (CL.1)	DB2 (CL.2)	CLASS 3 SUSPENSIONS	AVERAGE TIME	AV.EVENT	TIME/EVENT	HIGHLIGHTS	
ELAPSED TIME	0.001226	0.000510	LOCK/LATCH (DB2+IRLM)	0.000023	0.20	0.000116	#OCCURRENCES :	7841320
NONNESTED	0.001226	0.000510	IRLM LOCK+LATCH	0.000008	0.02	0.000347	#ALLIEDS :	0
STORED PROC	0.00000	0.00000	DB2 LATCH	0.000015	0.17	0.000085	#ALLIEDS DISTRIB:	0
UDF	0.00000	0.00000	SYNCHRON. I/O	0.000368	1.00	0.000368	#DBATS :	7841320
TRIGGER	0.00000	0.000000	DATABASE I/O	0.00000	0.00	0.000488	#DBATS DISTRIB. :	0
			READ CACHE HIT	0.00000	0.00	0.000488	#NO PROGRAM DATA:	0
CP CPU TIME	0.000103	0.000087	LOG WRITE I/O	0.000368	1.00	0.000368	#NORMAL TERMINAT:	7841320
AGENT	0.000103	0.000087	OTHER READ I/O	0.00000	0.00	N/C	#ROLLUP TRAN :	7841320
NONNESTED	0.000103	0.000087	OTHER WRTE I/O	0.00000	0.00	0.001171	#DDFRRSAF ROLLUP:	784132
STORED PRC	0.00000	0.000000	SER.TASK SWTCH	0.000001	0.00	0.072448	#ABNORMAL TERMIN:	0
UDF	0.000000	0.000000	UPDATE COMMIT	0.000000	0.00	N/C	#CP/X PARALLEL. :	0
TRIGGER	0.000000	0.000000	OPEN/CLOSE	0.000000	0.00	N/C	#UTIL PARALLEL. :	0
PAR. TASKS	0.00000	0.00000	SYSLGRNG REC	0.000000	0.00	N/C	#IO PARALLELISM :	0
SE CPU TIME	0.00000	0.00000	EXT/DEL/DEF OTHER SERVICE	0.000001	0.00	0.000288	#PCA RUP COUNT : #RUP AUTONOM. PR:	0

Figure 15: Omegamon XE for Db2 Performance accounting report – long.

About the authors

Tariq Hanif is a Senior Software Engineer in the z/OS I/O Supervisor Team. In z/OS, his focus areas of testing are Parallel Sysplex, Storage Replication Solutions and Disaster Recovery.

Brian Lee is a Software Developer for z/OS DFSMS Media Manager.

Dale Riedy is a Senior Technical Staff Member with z/OS I/O Supervisor Design and Development.

We would like to thank Nick Clayton, Distinguished Engineer - Enterprise Storage Development, for reviewing and providing comments on this article.

Rita Beisel contributed to the editorial review of this article.

1. <u>Faster Db2 Active Log Writes with Media Manager Parallel Write Support using</u> <u>zHyperLink</u>

رع

2. <u>z/OS MVS Initialization and Tuning Reference. Members of SYS1.PARMLIB.</u> <u>IECIOSxx (I/O related parameters).</u>

⇔

3. <u>z/OS MVS System Commands. MVS system commands reference. SETIOS</u> <u>command.</u>

<u>ب</u>

- 4. <u>Db2 Commands. -SET SYSPARM (Db2).</u> ←
- 5. <u>OMEGAMON for Db2 Performance Expert on z/OS 5.4.0. Accounting Report –</u> Long.