



Logical processor topology in HiperDispatch mode

Fabio Massimo Ottaviani – EPV Technologies

January 2017

1 Introduction

HiperDispatch's main goal is to maximize the probability for a workload to be re-dispatched to the same group of logical and physical processors, in order to allow it to reuse instructions and data previously loaded in the cache levels closer to the processor.

In order to decide the logical processor topology (mapping to the physical processors) HiperDispatch implements vertical polarization to split the LPAR logical processors in three pools:

- High polarity (high processor share); they will have a target share corresponding to 100% of a physical processor which will be pseudo-dedicated to each of them;
- Medium polarity (medium processor share); they will normally have a target share greater than 0% and less than 100% of a physical processor; these medium logical processors have the remainder of the LPAR's shares after the allocation of the logical processors with the high share; they will compete for the physical processors as before HiperDispatch;
- Low polarity (low processor share); they will receive a target share corresponding to 0% of a physical processor; these are considered discretionary logical processors which are not needed to allow the LPAR to fully utilize the physical processor resource associated with its weight; if there is not unused capacity left from other LPARs in the CEC they will be parked and not used.

By using the information provided in SMF 70 and 113 many consequences of the HiperDispatch activity can be analyzed. However, no information is provided in these records about the logical to physical processors topology.

This information is available in the SMF 99 subtype 14 (SMF 99-14) records and it can be easily collected in an SQL DB by EPV zParser customers.

However the IBM WLM team decided to make processor topology accessible to everyone and provided a free tool: the WLM Topology Report.

The tool provides an Excel spreadsheet that displays lots of interesting information such as:

- the association of logical processors to books, chips, drawers, and nodes,
- the vertical polarization of the processors (high, medium, low),
- the processor type (CPU and zIIP),
- the association to WLM affinity nodes,
- any topology changes.

In this paper, we will show what you have to do in order to install and use the WLM Topology Report. We will also show and discuss two real life examples.



2 Collecting the necessary SMF records

Information about HiperDispatch topology is provided in SMF 99-14. Records are written every 5 minutes or whenever a topology change occurs.

The most common reasons for a topology change are:

- Configuration changes
- Partition weight changes (which can be due to WLM soft capping).

A bit in the SMF99E_VCM_Flag1 field indicates if a topology change has occurred.

Records provide information by partition so to get a complete picture they have to be collected for all the partitions in the machine.

It's important to note that these records are always requested by IBM technical support when they are asked to investigate performance issues so it's highly recommended to collect them on a regular basis.

To collect SMF 99 subtype 14 records you have only to allow it in SMFPRMxx.

3 Downloading and installing the WLM Topology Report

The following pre-requisites apply:

- SMF 99-14 records belonging to a system running in HiperDispatch mode on a z10 or newer machine have to be available;
- MS Excel Version 2007, 2010 or 2013 has to be used.

The topology report tool can be downloaded from the IBM web site by using the following link:
http://www-03.ibm.com/systems/z/os/zos/features/wlm/WLM_Further_Info_Tools.html#Topology

The installation is quite simple; you have just to download and run the SetupTopologyReport.exe program. It will install the tool in the TopoReport folder in Windows Program Files and create the IBM RMF Performance Management shortcuts in the Windows start menu.

The topology report doesn't directly process SMF records. They have to be pre-processed on z/OS in order to produce a CSV file that is the needed input for the tool.

The following picture shows the content of the TopoReport folder.

Nome	Ultima modifica	Tipo	Dimensione
DocuTopo	14/01/2017 15:43	Cartella di file	
HostTopo	14/01/2017 15:43	Cartella di file	
Changes.txt	17/02/2016 16:36	File TXT	1 KB
TopoReport.xlsx	17/02/2016 16:36	Foglio con attivazi...	150 KB
uninst.exe	14/01/2017 15:43	Applicazione	35 KB

Figure 1



In the HostTopo folder, you will find two files:

- a load library, including the S99ERPTD program which converts SMF 99 subtype 14 records to CSV format;
- a JCL library, including the SAMPLE JCL to be used to run the program.

Nome	Ultima modifica	Tipo	Dimensione
TOPOREP.JCL.BIN	11/02/2015 13:16	BIN File	3 KB
TOPOREP.LOADLIB.BIN	06/12/2014 14:48	BIN File	12 KB

Figure 2

All the instructions to transfer these files to z/OS and use the tool can be obtained by clicking the Topo Report Help shortcut under the IBM RMF Performance Management entry in the Windows start menu.

4 Producing the CSV file

The SAMPLE JCL to be customized and run is provided below.

You have to:

- include your job card;
- set the name of the dataset including SMF 99-14 in SMFDSN;
- set the name of the output dataset including data in CSV format in RPTDSN;
- set the high level qualifier (hlq) for the TOPOREP.LOADLIB.

```
//      SET      SMFDSN=
//      SET      RPTDSN=
//JOBLIB      DD DSN=<hlq>.TOPOREP.LOADLIB,DISP=SHR
//*****{*}
//DECLARE      EXEC PGM=IEFBR14
//INPDS       DD DSN=&SMFDSN.,
//              DISP=SHR
//OUTDS       DD DSN=&RPTDSN.,
//              UNIT=SYSDA,
//              DISP=(MOD,DELETE),SPACE=(CYL,1)
//*****{*}
//READSMF     EXEC PGM=S99ERPTD
//SMFDATA     DD DISP=SHR,DSN=*.DECLARE.INPDS
//SYSPRINT    DD SYSOUT=*
//SYSUDUMP    DD SYSOUT=*
//REPORT      DD DISP=(NEW,CATLG),
//              SPACE=(CYL,(50,50),RLSE),
//              UNIT=SYSDA,
//              DSN=*.DECLARE.OUTDS
//USRPARMS   DD *
      ALLINTV
/*

```



If you want to load only the first record and the records where a topology change occurred, you have to comment out the ALLINTV parameter.

The resulting file has to be transferred in ASCII mode to the Windows system where the topology report tool is installed.

5 Understanding the report – Case 1

To produce the report you have to open the TopoReport.xlsxm spreadsheet. This is what you get.

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S															
1	Version:	1.001	WLM Topology Report																														
2	To begin a new Workbook: Create Copy Open a new CSV File and replace existing CSVREP: Open New CSV File Delete CSVREP: Delete CSVREP Worksheet																																
3																																	
4																																	
5																																	
6																																	
7	Interval Date Time System				Step 1: Copy Data to Main Worksheet								Step 2: Create Report																				
8	<input type="button" value="Clear Data"/>				<input type="button" value="Clear Report"/>								<input type="button" value="Sort Data"/>																				
9	<input type="button" value="Copy Data"/>				<input type="button" value="Create Report"/>								<input checked="" type="checkbox"/> Include System Name																				
10					<input type="checkbox"/> Clear Internals								<input type="button" value="Make Report"/>																				
11																																	
12																																	
13																																	
14																																	
15																																	
16																																	
17																																	
18																																	
19																																	
20																																	
21																																	
22																																	
23	Interval	Date	Time	System	Topo Chg	Rebuild	Flags	HonPrio Chg	WUQ Error	Speed Chg	CPs per Node	LPAR Share	# of CPs	CPU Index	CPU Type	Polarization	Aff. Node	Highest	NL1	Nestin													
24	Interval	Date	Time	System	Topo Chg	Rebuild	Flags	HonPrio Chg	WUQ Error	Speed Chg	CPs per Node	LPAR Share	# of CPs	CPU Index	CPU Type	Polarization	Aff. Node	Highest	NL1	NL2													

Figure 3

It's a good practice to create a new file to keep the original workbook clean by using the Create Copy button.

To load the CSV file, prepared as described in chapter 4, you have to click the Open New CSV File button.

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S															
1	Version:	1.001	WLM Topology Report																														
2	To begin a new Workbook: Create Copy Open a new CSV File and replace existing CSVREP: Open New CSV File Delete CSVREP: Delete CSVREP Worksheet																																
3																																	
4																																	
5	Interval Date Time System				Step 1: Copy Data to Main Worksheet								Step 2: Create Report																				
6	<input type="button" value="Clear Data"/>				<input type="button" value="Clear Report"/>								<input type="button" value="Sort Data"/>																				
7	<input type="button" value="Copy Data"/>				<input type="button" value="Create Report"/>								<input checked="" type="checkbox"/> Include System Name																				
8					<input type="checkbox"/> Clear Internals								<input type="button" value="Make Report"/>																				
9																																	
10																																	
11																																	
12																																	
13																																	
14																																	
15																																	
16																																	
17																																	
18																																	
19																																	
20																																	
21																																	
22																																	
23	Interval	Date	Time	System	Topo Chg	Rebuild	Flags	HonPrio Chg	WUQ Error	Speed Chg	CPs per Node	LPAR Share	# of CPs	CPU Index	CPU Type	Polarization	Aff. Node	Highest	NL1	Nestin													
24	Interval	Date	Time	System	Topo Chg	Rebuild	Flags	HonPrio Chg	WUQ Error	Speed Chg	CPs per Node	LPAR Share	# of CPs	CPU Index	CPU Type	Polarization	Aff. Node	Highest	NL1	NL2													

Figure 4

The window in the left section shows the records found. In this example, we collected only 5 records covering the time interval between 11:01:27 and 11:21:27 on February the 8th 2016. All the records belong to the TST1 system. In this system, all the SMF records are synchronized to 00. You can note that SMF 99-14 records are not synchronized instead.

Once you have selected the desired interval you have to click the Copy Data button to get the information provided in the lower section of the report.



A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
1	Version:	1,001																
2	To begin a new Workbook:		Create Copy		Open a new CSV File and replace existing CSVREP:			Open New CSV File		Delete CSVREP:		Delete CSVREP Worksheet						
3																		
4																		
5																		
6																		
7	Interval	Date	Time	System		Step 1: Copy Data to Main Worksheet		Step 2: Create Report										
8	1	2/8/2016	11:01:27	TST1														
9	2	2/8/2016	11:06:27	TST1														
10	3	2/8/2016	11:11:27	TST1														
11	4	2/8/2016	11:16:27	TST1														
12	5	2/8/2016	11:21:27	TST1														
13																		
14																		
15																		
16																		
17																		
18																		
19																		
20																		
21																		
22																		
23																		
24	Interval	Date	Time	System	Flags													
25	1	08/02/2016	11:01:27	TST1	Topo Chg	Rebuild	HonPrio Chg	WUQ Error	Speed Chg	CPs per Node	LPAR Share	# of CPs	CPU Index	CPU Type	Polarization	Aff. Node	Highest	Nesting
26	1	08/02/2016	11:01:27	TST1	No	No	No	No	No	6.625	15	0	0	High	1	2	1	1
27	1	08/02/2016	11:01:27	TST1	No	No	No	No	No	6.625	15	1	0	High	1	2	1	1
28	1	08/02/2016	11:01:27	TST1	No	No	No	No	No	6.625	15	2	0	High	1	2	1	1
29	1	08/02/2016	11:01:27	TST1	No	No	No	No	No	6.625	15	3	0	High	1	2	1	1
30	1	08/02/2016	11:01:27	TST1	No	No	No	No	No	6.625	15	5	0	Med	1	2	5	1
31	1	08/02/2016	11:01:27	TST1	No	No	No	No	No	6.625	15	6	0	Med	1	2	5	1
32	1	08/02/2016	11:01:27	TST1	No	No	No	No	No	6.625	15	7	0	Low	1	2	5	1
33	1	08/02/2016	11:01:27	TST1	No	No	No	No	No	6.625	15	8	5	High	2	2	4	1

Figure 5

The Flags columns tell you if specific events happened in the selected interval:

- Topo Chg, topology has changed;
- Rebuild, affinity nodes have been rebuilt;
- HonPrio Chg, the honour priority parameter has been changed;
- WUQ Error, an error occurred on a dispatcher queue;
- Speed Chg, the processor speed changed.

The most relevant information provided in the next columns are the logical processor type (0=CPU, 5=zIIP), its polarization and the WLM affinity node they are assigned to.

Scrolling to the right, in the Nesting Levels columns you will get information about the processor cache architecture of the machine.

TST1 runs on a zEC12 so there are only two nesting levels: NL1 (Chip) and NL2 (Book).

For each logical processor, you know which Chip and Book it is dispatched onto. At the moment, no information about the physical processor is provided.

You may see that the first 5 logical processors are CPUs with High polarization and they are all assigned to Chip 1 and Book 1.

I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB
1	WLM Topology Report																		
2	Open a new CSV File and replace existing CSVREP:			Delete CSVREP: Delete CSVREP Worksheet															
3	Worksheet					Step 2: Create Report													
4																			
5																			
6																			
7																			
8																			
9																			
10																			
11																			
12																			
13																			
14																			
15																			
16																			
17																			
18																			
19																			
20																			
21																			
22																			
23																			
24	Speed Chg	CPs per Node	LPAR Share	# of CPs	CPU Index	CPU Type	Polarization	Aff. Node	Highest	NL1	NL2	NL3	NL4	NL5	New STSI Format	Version	SMF Subversion	Description	
25	No	6.625	15	0	0	High	1	2	1	1	1	0	0	0	27	26	NL1=Chip	NL2=Book	
26	No	6.625	15	1	0	High	1	2	1	1	0	0	0	0	27	26	NL1=Chip	NL2=Book	
27	No	6.625	15	2	0	High	1	2	1	1	0	0	0	0	27	26	NL1=Chip	NL2=Book	
28	No	6.625	15	3	0	High	1	2	1	1	0	0	0	0	27	26	NL1=Chip	NL2=Book	
29	No	6.625	15	4	0	High	1	2	1	1	0	0	0	0	27	26	NL1=Chip	NL2=Book	
30	No	6.625	15	5	0	Med	1	2	5	1	0	0	0	0	27	26	NL1=Chip	NL2=Book	
31	No	6.625	15	6	0	Med	1	2	5	1	0	0	0	0	27	26	NL1=Chip	NL2=Book	
32	No	6.625	15	7	0	Low	1	2	5	1	0	0	0	0	27	26	NL1=Chip	NL2=Book	
33	No	6.625	15	8	5	High	2	2	4	1	0	0	0	0	27	26	NL1=Chip	NL2=Book	

Figure 6



By clicking the Make Report button, you will get a clearer picture of the logical to physical processor topology.

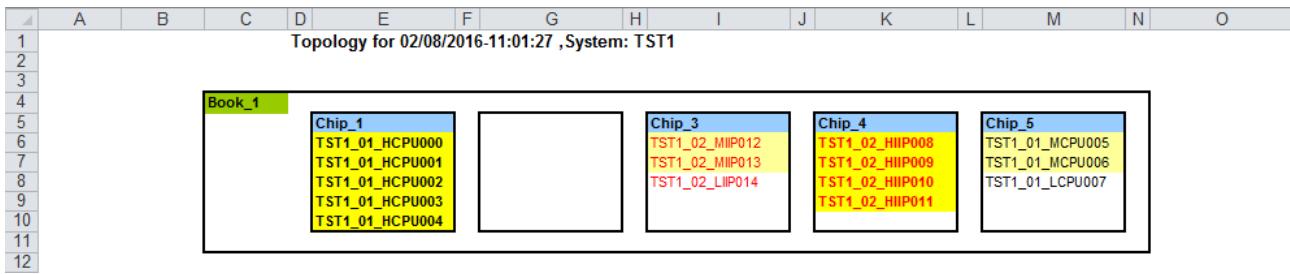


Figure 7

“The logical processors are depicted with the following information: SSSS_NN_Vtttnnn which means:

- SSSS = SMF system id
- NN = WLM affinity node number
- V = polarization = {H,M,L}
- ttt = processor type = {CPU,IIP,AAP}
- nn = processor number

High processors are also depicted in yellow, medium processors in light yellow and zIIPs with red color.¹

You can note that:

- all the TST1 logical processors are dispatched to Book 1;
- WLM affinity node 1 includes all the CPUs: the 5 CPUs with High polarization are dispatched to Chip 1 while the other CPUs (2 Medium and 1 Low) are dispatched to Chip 5;
- WLM affinity node 2 includes all the zIIPs: the 4 zIIPs with High polarization are dispatched to Chip 4 while the other zIIPs (2 Medium and 1 Low) are dispatched to Chip 3.

6 Understanding the report – Case 2

In this second example, we refer to systems running on a big z13 machine.

Figure 7 shows the topology for the SYS1 system. You can note that both the CPUs and the zIIPs are split on 2 Drawers and 3 Nodes (even if all the High are in Drawer 3 and most of them in Drawer 3 - Node 1).

¹ From the WLM Topology Report help.



	A	B	C	D	E	F	G	H	I	J	K
1	Topology for 01/05/2017-00:02:39 , System: SYS1										
Drawer_3											
3	Node_1		Chip_1	SYS1_02_HCPU023 SYS1_02_HCPU024 SYS1_04_HCPU025 SYS1_04_HCPU026 SYS1_07_HCPU027 SYS1_07_HCPU028	Chip_2	SYS1_03_HCPU005 SYS1_03_HCPU006 SYS1_08_HCPU007 SYS1_08_HCPU008 SYS1_09_HCPU009 SYS1_09_HCPU022	Chip_3	SYS1_10_HCPU002 SYS1_10_HCPU003 SYS1_10_HCPU004 SYS1_05_HIIP013			
4											
5											
6											
7											
8											
9											
10											
11											
12											
13											
14	Node_2		Chip_1	SYS1_01_MCPU029 SYS1_01_LCPU030	Chip_2	SYS1_01_HCPU000 SYS1_01_HCPU001 SYS1_05_HIIP014	Chip_3	SYS1_05_HIIP010 SYS1_05_HIIP011			
15											
16											
17											
18											
19											
20											
21											
22											
23											
24	Drawer_4		Node_2		Chip_2	SYS1_05_HIIP015 SYS1_05_HIIP016 SYS1_05_HIIP017	Chip_3	SYS1_05_HIIP018			
25											
26											
27											
28											
29											
30											
31											
32											
33											

Figure 7

The tool doesn't allow you to select more than one system at a time so you have to do more reports to get a complete picture.

Figure 8 shows the topology for the SYS2 system. You can note that all the CPUs and zIIPs are assigned to Drawer 2. CPUs are split between the two Nodes while most of the zIIPs are in Drawer 2 - Node 1.

	A	B	C	D	E	F	G	H	I	J	
1	Drawer_2										
Node_1											
2				Chip_1	SYS2_04_HIIP010 SYS2_04_HIIP011 SYS2_05_HIIP012 SYS2_05_HIIP013 SYS2_05_HIIP015 SYS2_10_HCPU026	Chip_2	SYS2_01_HCPU022 SYS2_01_HCPU023 SYS2_03_HCPU024 SYS2_03_HCPU025 SYS2_06_HCPU027 SYS2_06_HCPU028	Chip_3	SYS2_07_HCPU006 SYS2_07_HCPU007 SYS2_08_HCPU008 SYS2_08_HCPU009		
3											
4											
5											
6											
7											
8											
9											
10											
11											
12	Node_2			Chip_1	SYS2_02_HCPU000 SYS2_02_HCPU001 SYS2_09_HCPU002 SYS2_09_HCPU003 SYS2_10_HCPU004 SYS2_10_HCPU005	Chip_2	SYS2_04_HIIP014 SYS2_02_MCPU029 SYS2_02_LCPU030				
13											
14											
15											
16											
17											
18											
19											
20											

Figure 8

We also produced the SYS3 report and we discovered that it uses the same drawers, nodes and, in some cases, the same chips as SYS1. So we manually integrated SYS1 and SYS3 reports and highlighted the shared chips with a green border.



A	B	C	D	E	F	G	H	I	J
23	Drawer_3								
24		Node_1							
25			Chip_1		Chip_2		Chip_3		
26			SYS1_02_HCPU023		SYS1_03_HCPU005		SYS1_10_HCPU002		
27			SYS1_02_HCPU024		SYS1_03_HCPU006		SYS1_10_HCPU003		
28			SYS1_04_HCPU025		SYS1_08_HCPU007		SYS1_10_HCPU004		
29			SYS1_04_HCPU026		SYS1_08_HCPU008		SYS1_05_HIIP013		
30			SYS1_07_HCPU027		SYS1_09_HCPU009		SYS1_09_HCPU022		
31			SYS1_07_HCPU028						
32									
33									
34		Node_2			Chip_1		Chip_3		
35			SYS1_01_MCPU029		SYS1_01_HCPU000		SYS1_05_HIIP010		
36			SYS1_01_LCPU030		SYS1_01_HCPU001		SYS1_05_HIIP011		
37					SYS1_05_HIIP014		SYS1_05_HIIP012		
38							SYS3_10_MCPU029		
39							SYS3_10_LCPU030		
40									
41									
42									

Figure 9

A	B	C	D	E	F	G	H	I	J
45	Drawer_4								
46		Node_1			Chip_1		Chip_3		
47			SYS3_01_HCPU023		SYS3_02_HCPU005		SYS3_03_HCPU000		
48			SYS3_01_HCPU024		SYS3_02_HCPU006		SYS3_03_HCPU001		
49			SYS3_04_HCPU025		SYS3_08_HCPU007		SYS3_10_HCPU002		
50			SYS3_04_HCPU026		SYS3_08_HCPU008		SYS3_10_HCPU003		
51			SYS3_07_HCPU027		SYS3_09_HCPU009		SYS3_03_HCPU004		
52			SYS3_07_HCPU028		SYS3_09_HCPU022				
53									
54									
55									
56		Node_2			Chip_1		Chip_3		
57			SYS3_05_HIIP010		SYS1_05_MIIIP015		SYS1_05_LIIP018		
58			SYS3_05_HIIP011		SYS1_05_LIIP016		SYS1_05_LIIP018		
59			SYS3_06_HIIP012		SYS1_05_LIIP017		SYS3_05_MIIIP016		
60			SYS3_06_HIIP013		SYS2_05_LIIP016		SYS3_05_MIIIP017		
61			SYS3_06_MIIIP014		SYS2_05_LIIP017		SYS3_05_MIIIP016		
62							SYS3_05_MIIIP017		
63									
64									
65									
66									

Figure 10

More systems run on this machine. To get the complete picture we should create a report for each system and integrate them manually.

7 Conclusions

The tool can be very useful to understand the processor placement and how it changes when topology and other changes occur.

Most relevant current limitations are:

- SMF records are not synchronized with other SMF and RMF records;
- No information about logical to a specific physical processor is provided;
- Only one system at a time can be selected.