



CPU Latent Demand and Specialized Processors

Fabio Massimo Ottaviani
EPV Technologies

Performance analysts know that, if the CPU utilisation is at 100% for long periods of time, there is a strong probability that the system is saturated and a further resource demand, hidden in the CPU queue, needs to be estimated. This demand is generally known as “CPU latent demand”.

Most of the techniques used to estimate latent demand in the last years have been based on the CPU ready queue. Unfortunately the availability of the new specialized processors (zAAP and zIIP) changed the meaning of this metric: it now provides the number of AS running on or waiting for GCP, zAAP or zIIP (not only for GCP).

This paper will discuss latent demand in the new zAAP and zIIP world suggesting an approach to understand if latent demand is in the system and to which kind of processor it belongs to.

INDEX

1	THE CPU READY QUEUE	2
2	IS THERE CPU LATENT DEMAND ?	2
3	GCPS AND ZAAPS SCENARIOS	4
4	CONCLUSIONS	4



1 The CPU ready queue

One of the most interesting metrics available in SMF 70 and in the RMF CPU report is the number of IN&READY address spaces. It shows the number of address spaces running on a processor plus the number of address spaces ready to be executed.

When the number of IN&READY address spaces is much bigger than the number of available processors it means that a long queue in front of each processor exists.

The first consequence of this situation is that some activities in the system (generally the ones at the lowest importance) will be slowed down.

A second issue to consider is that the measured CPU busy (close to 100% for a very long period of time) is not a good estimate of the real CPU needs anymore.

What is missing is the “CPU Latent Demand” that is the CPU needed by the address spaces in the IN&READY queue to run at the appropriate speed. This latent demand has to be considered in capacity planning studies in order to avoid a severe under estimate of the CPU capacity.

Obviously a certain level of queuing has to be accepted if you want to exploit the full machine capacity so a threshold is needed to understand when the queuing is excessive and latent demand has to be estimated.

A generally accepted ROT says that if the IN&READY value is more than 2 times the number of available processors then latent demand is consistent and needs to be evaluated.

Unfortunately the availability of the new specialized processors (zAAP and zIIP) changed the meaning of the IN&READY metric: it now provides the number of AS running on or waiting for any of the different processor types (GCP, zAAP or zIIP).

2 Is there CPU Latent Demand ?

In order to make things easier, only GCP and zAAP will be discussed in the following. However the same concepts presented for zAAP applies to zIIP as well.

Figure 1 shows a snapshot of the RMF CPU activity report. The number of address spaces measured in the IN-READY queue are on average 13,5 (minimum is 2, maximum is 67).

The number of online processors is 12,8 (decimal values are normal if IRD or CoD are active) so the ratio between the number of queued address spaces and the number of processors is approximately 1.



3 GCPs and zAAPs scenarios

Depending on the load of GCPs and zAAPs, on the number of CPs of each type and on the IFAHONORPRIORITY parameter setting, more scenarios have to be evaluated in order to understand if latent demand is present in the system and to which processor type it belongs¹.

a) Light GCP load and light zAAP load (IFAHONORPRIORITY=YES/NO)

No latent demand.

b) Light GCP load and heavy zAAP load (IFAHONORPRIORITY=YES)

GCP will help zAAPs to serve the workloads. The IN-READY queue values should be below the threshold and no consistent latent demand should be in the system. However you have to check if you are using expensive GCP resources instead of the cheaper zAAPs looking at the amount of zAAP eligible work running on GCPs. This work will become latent demand on zAAPs if the IFAHONORPRIORITY parameter would be set to NO.

c) Light GCP load and heavy zAAP load (IFAHONORPRIORITY=NO)

GCP will not help zAAPs to serve the workloads. Checking the IN-READY queue values is not enough to understand if zAAPs latent demand is present in the system. (see Scenario 1 in Chapter 2).

d) Heavy GCP load and light zAAP load (IFAHONORPRIORITY=YES/NO)

Checking the IN-READY queue values is not enough to understand if GCPs latent demand is present in the system. (see Scenario 2 in Chapter 2).

e) Heavy GCP load and heavy zAAP load (IFAHONORPRIORITY=YES)

GCP will help zAAPs to serve the workloads. If the IN-READY queue is below the threshold you could assume that no significant latent demand is present in the system but you have to check if you are using expensive GCP resources instead of the cheaper zAAPs looking at the amount of zAAP eligible work running on GCPs. If the IN-READY queue is above the threshold latent demand could be on GCPs, on zAAPs or on both of them.

f) Heavy GCP load and heavy zAAP load (IFAHONORPRIORITY=NO)

GCP will not help zAAPs to serve the workloads. Checking the IN-READY queue values is not enough to understand if GCPs or zAAPs latent demand is present in the system. (similar to Scenario 1 in Chapter 2).

4 Conclusions

The availability of specialized processors makes it much more complex to estimate CPU latent demand. The main problem is that the number of IN-READY address spaces is not provided per CP type. The old methods used to estimate latent demand have to be deeply revised and empowered using a different approach taking into consideration other elements and metrics.

In this paper a first step is suggested to understand if latent demand is in the system and to which processor type it belongs to.

¹ The addition of zIIPs will substantially increase the number of different scenarios to consider.