# WLM Enhancements for z14 and z/OS 2.3

Horst Sinram

STSM, z/OS Workload and Capacity Management
sinram@de.ibm.com

02 Dec 2017

# Content

❑ New Capping Options

❑ Reporting Enhancements for CICS, IMS, and Mobile Workloads

❑ Container Pricing*

❑ More granular resource controls

❑ IBM z14 exploitation

❑ z/OS 2.3 Preview

IBM

# New Capping Options

| z/OS release / Function | V2.2 | V2.1 |
|---|---|---|
| z13 GA2 LPAR Absolute group capping | OA47752 | OA47752 |
| Absolute MSU capping | OA49201 | OA49201 |

❑ **LPAR absolute group capping**
 ➢ Feature of PR/SM as of z13 GA2, and z13s

 ➢ Like LPAR absolute capping but for a group of LPARs

 ➢ Specified on the HMC as number of processors with 2 decimal places (like 3.75). All processor classes supported.
 ➢ Recognized by WLM as possible limit to the LPAR capacity

❑ **WLM Absolute MSU capping**
 ➢ Function of WLM provided by APAR OA49201. Requires zEC12 GA2 or later.

 ➢ Similar to WLM defined capacity or group capacity but LPAR will always be capped
  ○ Independent of 4 hour rolling average consumption.
  ○ General purpose processor
 ➢ Specified in IEAOPTxx.
  ○ Limit is the LPAR defined capacity or group capacity specified on the HMC **in MSU**.

IBM Z

# Using Absolute MSU Capping

| IEAOPTxx ABSMSUCAPPING= | |
|---|---|
| <u>NO</u> | Defined capacity limits and group capacity limits should be enforced only while the long term four hour rolling average consumption exceeds the respective limit (existing and usually desired behavior). |
| YES | Defined capacity limits and group capacity limit should be enforced **permanently, independently of the long term four hour rolling average consumption**. Becomes effective on zEC12 GA2 or later. |

❑ ABSMSUCAPPING=Yes limits LPAR consumption to a certain MSU number at all times.

➢ The system loses the flexibility of consuming above the defined capacity limit while the four hour rolling average is below the limit.

❑ Limit remains stable even when CEC configuration changes, e.g. through On/Off CoD or CBU activations or deactivations.

➢ Absolute MSU capping is an effective means to permanently limit the consumption of an LPAR to a specific MSU figure at all times

○ Including times when the *four-hour rolling average* does not exceed the defined limit.

# Using Absolute MSU Capping with Group Capacity

❑ When used with an LPAR capacity group:
  ➢ Limit on behalf of the group entitlement will always be enforced
    ○ Regardless of the *four-hour rolling group average* consumption.

  ➢ As with ABSMSUCAPPING=NO, an LPAR is allowed to take benefit of the unused group capacity
    ○ Unless the LPAR is also capped via other LPAR limits.

  ➢ All members of a capacity group that use ABSMSUCAPPING=YES will permanently enforce the limit on behalf of the capacity group.

  ➢ All members of a capacity group that do *not* use ABSMSUCAPPING=YES will be capped while the group *four-hour rolling group average* consumption is greater or equal to the group limit

# Reporting Enhancements for CICS, IMS, Mobile and Cloud Workloads

- ❑ Mobile Workload Pricing (MWP) is an IBM Software Pricing option, announced in May 2014
- ❑ Workload Pricing for Cloud (zWPC) is an IBM Software Pricing option, announced in July 2016

- ❑ For eligible software both can reduce the cost of transactions that originate from mobile devices or new public cloud workloads
  - ➤ MWP and zWPC can mitigate the impact of such workloads on sub-capacity license charges, specifically in the cases where higher mobile or cloud transaction volumes may cause a spike in machine utilization
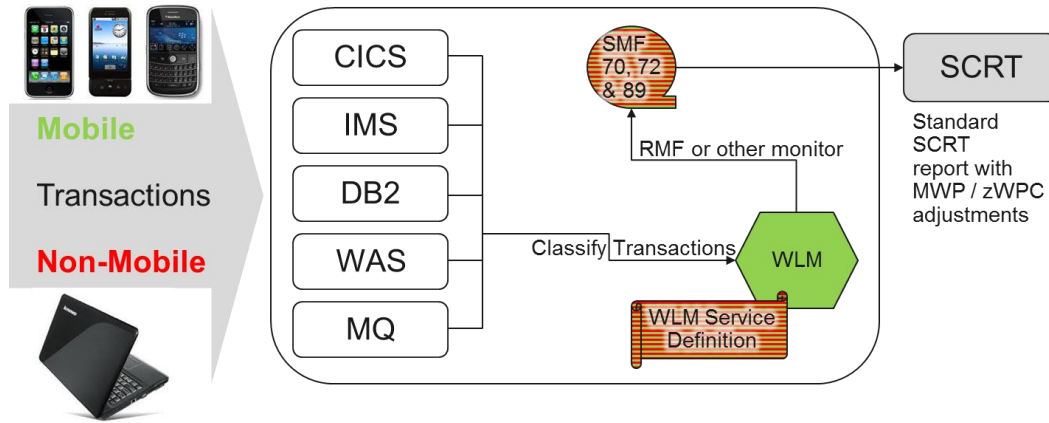
- ❑ Reporting enhancements for CICS, IMS, Mobile, and Cloud Workloads introduce WLM enhancements that can simplify the identifying and reporting of the mobile- or cloud-sourced transactions and their processor consumption

# Reporting Enhancements for Mobile Workloads

| z/OS release / Function | z/OS V2.1 | z/OS V2.2 | Other |
|---|---|---|---|
| WLM: Reporting Enhancements for Mobile Workloads | *OA47042* | *OA47042* | |
| RMF Reporting Enhancements for Mobile Workloads | *OA48466* | *OA48466* | |
| *z/OSMF* Reporting Enhancements for Mobile Workloads | *PI47638* | *PI47638* | |
| *CICS TS* Reporting Enhancements for Mobile Workloads | | | *CICS 5.3* |
| *IMS TM* Reporting Enhancements for Mobile Workloads | | | *IMS 14 PI46933 (available) PI51948 (1H2016\*)* |
| *SCRT and Billing System Support* | | | *SCRT 23.13.0 (mobile, available) SCRT 24.10.0 (cloud, available)* |

IBM Z

# Identifying Mobile and Cloud Workloads



- ❑ MWP and zWPC offer a discount on MSUs consumed by transactions that originated from a mobile device or new public cloud workloads

- ❑ To take advantage of this discount, you need a process, agreed upon by you and IBM, **to identify (tag and track)** mobile- or cloud-sourced transactions and report on their consumption

- ❑ NEW Identify mobile or cloud transaction via a transaction level attribute in the WLM service definition
  - ➤ Processor consumption data aggregated by WLM
  - ➤ Reporting integrated into standard performance monitors (RMF) and low volume SMF records
  - ➤ Applicable to wide range of workloads, including enclave work and CICS/IMS work

# Identifying Mobile and Cloud Workloads …

❑ In your WLM classification rules, classify transactions as **MOBILE, CATEGORYA or CATEGORYB**

❑ The assigned attribute is independent from the assigned service and report class
  ➤ Eliminates the need for using new dedicated classes for mobile or cloud workload reporting

❑ The assigned attribute is transparent to subsystems

❑ WLM tracks and reports the total and the **MOBILE, CATEGORYA and CATEGORYB CPU consumption** for all service and report classes
  ○ With exploiting levels of CICS and IMS, CPU consumption data is also available for CICS and IMS transaction service and report classes that previously did not report any CPU consumption data
  ○ Subsystems using independent enclaves can participate transparently; only the classification rules need to be updated.

❑ WLM also aggregates and reports the **system-wide MOBILE, CATEGORYA and CATEGORYB consumption** data
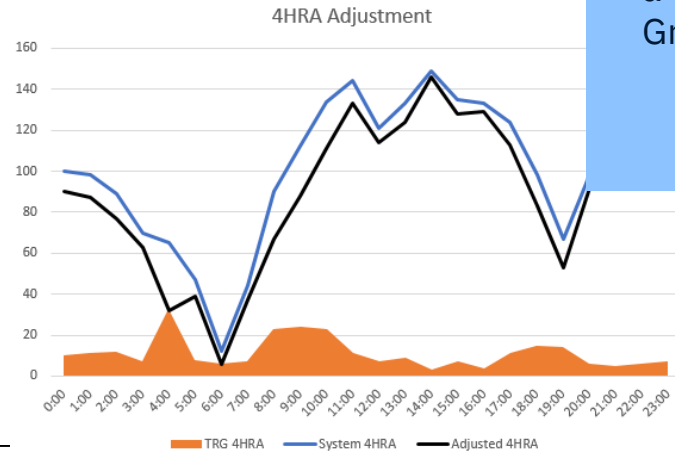
# Container Pricing for IBM Z

- **Announcement** excerpt:
  IBM is introducing Container Pricing for IBM Z for qualified solutions running on IBM z13 and z14 servers. Container Pricing will provide simplified software pricing for qualified solutions, combining flexible deployment options with competitive economics that are directly relevant to those solutions.

- Container Pricing can scale from collocated solutions within existing LPARs, through to separate LPARs, up to multiple-LPAR solutions, without directly impacting the cost of unrelated workloads. Additionally, Container Pricing will simplify pricing and billing on the IBM Z platform, by superseding a number of existing price offerings and by fully automating the billing process.

- IBM initially announces three solutions that will be enabled with Container Pricing:
  - The New Application Solution will provide a highly competitive stand-alone priced offering for new z/OS applications, such as CICS® TS or WebSphere applications. The New Application Solution is the strategic replacement for the current zWPC and IWP priced offerings.
  - The Application Development and Test Solution will provide highly competitive stand-alone pricing for z/OS based development and test workloads. Modern DevOps tooling can be optionally added at uniquely discounted prices.
  - The Payments Solution will provide a "per payment" pricing option for IBM Financial Transaction Manager for z/OS deployments. This new offering directly ties operational cost to business value by basing the price on the number of payments processed, rather than capacity used to process them.

- Container Pricing for IBM Z is planned to be available by year end 2017 and enabled in z/OS V2.2 and z/OS V2.3 with the PTFs for APARs associated with fix category IBM.Function.PricingInfrastructure. z/OS will enhance both the Workload Manager capability of z/OS (z/OS WLM) and the Sub-Capacity Reporting Tool (SCRT) to support Container Pricing. This includes:
  - The introduction of a new Tenant Resource Group capability within z/OS WLM to allow the metering and optional capping of workloads, along with the ability to map those workloads directly to Container Pricing.
  - Enhancements to SCRT to capture eligible Container Pricing workloads, allowing for the billing of those solutions independently of traditional Sub-Capacity pricing.

- For more information, see Whitepaper or https://www-03.ibm.com/systems/z/resources/swprice/container.html

# Motivation for the z/OS service definition objects?

Cloud workload paradigm asks for new ways of metering workloads in multi-tenant environments

IBM Z business asks for an infrastructure to support novel pricing options

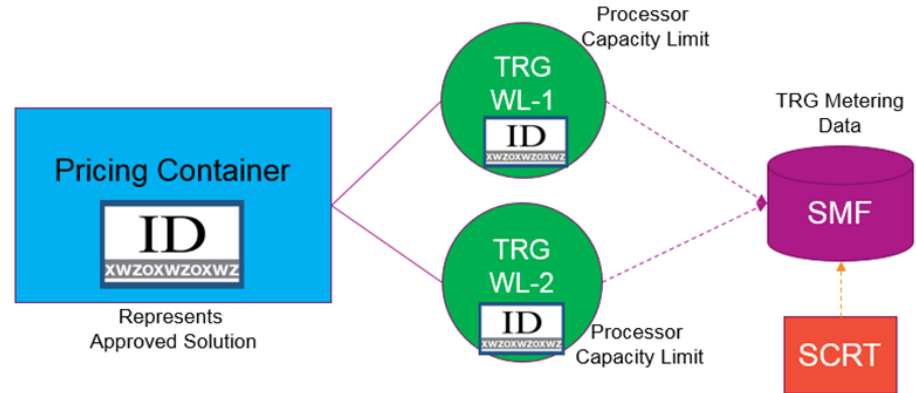**4HRA Adjustment**



- A Tenant Report Class (TRC) is similar to a WLM Report Class. TRCs are assigned through WLM classification and are always associated with a Tenant Resource Group.

- A Tenant Resource Group (TRG) is somewhat similar to a WLM Resource Group and can be associated with tenants or solutions.

- TRGs aggregate consumption data and can optionally be used to apply consumption limits.

# Container Pricing workflow overview

- For a eligible solution, IBM provide customers with a "Solution ID" (key)

- In their WLM service definition customers
  - Define one or more TRGs and paste the Solution ID into the definition
    - Dummy solution IDs for test/education will be documented
  - Define one or more TRCs associated with TRG
  - Change or add classification rules to classify eligible work (only!) and assign service class and TRCs
  - Install and activate WLM service definition
- Monitors query WLM and write new SMF70 data sections for TRGs
  - In addition, the TRC and TRG data will be reported via the existing report class and resource group mechanism in the WLMGL report

- SCRT consumes SMF70 and SMF89 data for billing
  - Verifies solution ID, applies pricing rules

# Support for IBM Container Pricing for IBM Z

| z/OS release / Function | z/OS V2.3 | z/OS V2.2 | z/OS V2.1 |
|---|---|---|---|
| WLM | OA52312 | OA52312 | Coexistence only: OA52312 |
| RMF | OA52694 | OA52694 | |
| z/OSMF WLM | PI89361 | PI89361 | |
| z/OSMF RMF | PI89935 | PI89935 | |
| SMF | OA53033 | OA53033 | |
| SDSF | PI82528 | PI82528 | |
| SCRT and Billing System Support | OA53047 | OA53047 | |

Container Pricing FIXCAT category: Keyword:
Content collection (Knowledge center)
IBM.Function.PricingInfrastructure
PRICINGINFR/K

# TRG and TRC Definition

- The WLM Administrative Application Level is increased to 32.

- Tenant Resource Groups and Tenant Report Classes can be defined via new menu items.

- Specification of these new objects will increase the functionality level of the service definition to 32.

```
Functionality LEVEL032       Definition Menu        WLM Appl LEVEL032


---------------------------------------------------------------


Definition data set  . . : 'WLM.DEMO.SRVDEF.XML'
Definition name  . . . . . PROD01    (Required)
Description  . . . . . . . Production service definition


Select one of the following options.

__      1.  Policies               12. Tenant Resource Groups
        2.  Workloads              13. Tenant Report Classes
        3.  Resource Groups
        4.  Service Classes
        5.  Classification Groups
        6.  Classification Rules
        7.  Report Classes
        8.  Service Coefficients/Options
        9.  Application Environments
        10. Scheduling Environments
        11. Guest Platform Mgmt Provider
```

# Tenant Resource Group (TRG) Definition

- **The TRG name is mandatory (8 char)**

- **Description, Tenant ID, Tenant Name are optional and are expected to be used in a z/OS cloud context**

- **For qualified offerings, a 64 char Solution ID needs to be provided.**

```
               Create a Tenant Resource Group
Enter or change the following information:


Tenant Resource Group Name   TRGDEM01
Description . . . . . . . . Sample TRG
Tenant ID . . . . . . . . . _____
Tenant Name . . . . . . . . Solution Newapp
Solution ID . . . . . . . .
     Z194E15-F1078F4-CEBBF9F075-099853BF-60EF-4A05-A0D1-EF925B-992C90

Define Capacity: __  1.  In Service Units (Sysplex Scope)
                     2.  As Percentage of the LPAR share (System Scope)
                     3.  As a Number of CPs times 100 (System Scope)
                     4.  In accounted workload MSU (Sysplex Scope)
Maximum Capacity . . . . . . . . . . . _____
Include Specialty Processor Consumption NO      (YES or NO)
```

- Exactly enter (paste) the IBM provided Solution ID string.
  – WLM user interfaces perform only sanity check the ID. Solution IDs failing that check are rejected.
  – Attributes encoded into the ID may change (or not) how the system processes the work.
  – The Solution ID is acted upon during SCRT processing.
  – Multiple TRGs may specify same Solution ID

- Optionally, a consumption limit can be specified. TRG capacity limits should not be specified unless there is a need to limit processor consumption.
  – The "Include Specialty Processor Consumption" switch indicates whether the combined CP and specialty processor consumption determines the cap limit.

- Unlike standard resource groups
  – there is no minimum consumption limit, no memory limit

# Tenant Report Class (TRC) Definition

- The TRC name is mandatory (8 char)
  - Name must be unique (also across report classes)
  - In total, up to 2047 Report Classes and Tenant Report Classes can be defined

```
         Create a Tenant Report Class

---------------------------------------------------

Enter or change the following information:

Tenant Report Class Name . . . T_CDC      (Required)
Description  . . . . . . . . . TRC including CDC region

Tenant Resource Group Name . . TRGDEM01   (Required; name or ?)
```

- The TRG name is required, i.e. any TRC must be associated with a TRG

- Monitoring interfaces and monitors report on TRC as on standard report classes

# TRC related classification rules considerations

■ Reporting Attributes MOBILE, CATEGORYA&B must not be used with a Tenant Report Classes in the same classification rule

■ A Tenant Report Class must not be used with a service class that is associated with a Resource Group in the same classification rule

■ If a Tenant Report Class is used in classification rules that assign different service classes, the Tenant Report Class might become heterogeneous

  o This means that work may run in different service classes but reported altogether in this one Tenant Report Class
  o WLM validation issues a warning panels and messages, such as
  ```
  IWMAM916W Tenant Report Class T_CDC might become heterogeneous by combining
  work running in service classes VEL80 VEL50
  ```
  o Strong recommendation is to use only homogeneous TRCs:
  create different Tenant Report Classes for each service class, and connect them all to the same Tenant Resource Group

# Comparison of TRG and RG Capping Types

All RGs and TRGs are Sysplex-wide defined.
The limit may be evaluated either on the Sysplex level or on each system.
For **all** (T)RG types only captured TCB and SRB times are counted towards the limit.
The limit is enforced based on a one minute average (i.e., no 4HRA).

| Type 1 | Type 2 | Type 3 | Type 4 |
|---|---|---|---|
| Raw CPU+SRB service units ("Raw" meaning that Service Definition Coefficients are not applied). | Percent of CP LPAR share (even if specialty processor consumption included). May exceed 100%. | Percent of one CP processor (even if specialty processor consumption included). | Processor consumption expressed in "accounted workload MSU" – see RG and TRG "MSU" limits. |
| Limit applies to Sysplex. | Limit applies to each System. | Limit applies to each system. | Limit applies to Sysplex |

Up to 32 RGs **plus 32 TRGs** may be defined.

# TRG Capping

TRG capping is based on Resource Group capping and inherits its characteristics.

- Time is divided into 256 "slices". In any slice the whole (T)RG can be set
  -dispatchable (called awake slice)
  -non-dispatchable (cap slice)

- Regardless "`Include Specialty Processor Consumption`" setting, work will be not dispatchable during a cap slice on any processor types.

- The cap pattern is adjusted every 10 sec based on the average of the last minute

- (T)RG may be comprised of work at different priority. Cap pattern applies to entire (T)RG, i.e. during awake slices higher priority work will be dispatched ahead of lower priority work.

- (T)RG consumption will vary based on demand, mix of dispatch priorities, number of dispatchable units and number of processors.
  - The system will attempt to over-cap the work, i.e., the consumption will be throttled to remain below the limit. Depending on the characteristics of the work this may not always be possible.
  - Usually, consumption levels in within (few) minutes

- Very latency sensitive work not a good candidate for capping. Multiple TRGs with same Solution ID may be used when needed.

Sample cap pattern showing work that is capped 50% of time.



**Tenant Resource Group Overview**
Tenant Resource Group: TRGTLL1 , Type2: 0-30% of LPAR

# RG and TRG limits
## …especially "accounted MSU" limits

- Background: Technical and pricing related performance/capacity numbers are based on different views. This remains unchanged.
  - The adjustment factor for service units (technical view) is based on the logical configuration (number of CP s online to the LPAR) on the respective CPC. Refer to Processor version codes and SRM constants
  - The adjustment factor for pricing purposes (MSU) is based on the physical configuration (CPC model capacity rating). Refer to Large Systems Performance Reference for IBM Z

- Every 10 sec, WLM converts type 2, 3, and 4 (T)RG limits into a SU/sec service rate based on current configuration.
  - Therefore, a type 4 (MSU) limit will be converted using the CEC and LPAR adjustment factors.
  - A type 4 limit is intended to simplify the specification of a limit expressed in MSU, **but…**

- **It cannot be expected that RG and TRG MSU limits will closely match the resulting LPAR MSU consumption or $HRA TRG consumption:**
  - **The (T)RG limit applies only to the accounted (captured) TCB and SRB times. System management time (uncaptured time) is not included.**
  - **The limit is not a 4 hour rolling average.**

# When to use resource groups or tenant resource groups

- **Only or preferentially use standard resource groups...**

  - When a resource group minimum or memory limit is required

  - When it is required to use WLM managed initiators as part for the RG

- **Only or preferentially use tenant resource groups...**

  o For authorized pricing container solutions

  o For all functions when the group just serves aggregation (i.e. no limits)

# Comparison of Interim solution (CategoryA/B) vs. TRG solution

## CategoryA/B solution

- **Granularity implied by WLM classification rule**

- **4HRA aggregated and propagated into SMF70. System overhead included (apportioned)**

- **Metering only**

- **Up to 2 categories**

- **CategoryA/B may be phased out in the future. Continued to be used for mobile workload reporting in the immediate future.**

## TRG solution

- **Granularity implied by WLM classification rule. Migration from interim solution very straight- forward.**

- **4HRA aggregated and propagated into SMF70. System overhead included (apportioned)**

- **Metering and (optional) capping**

- **Up to 32 TRGs may be defined.**
  - **IBM provided solution IDs to map to T&Cs**
  - **SMF89 (SCRT) correlation**
- **Strategic solution**

# Comparison of capping types

| Type of capping | Scope | Specification unit | Proc types | Suitable to isolate LPAR(s), or to enforce hard consumption limit | Manage-able by CPM** | Control point |
|---|---|---|---|---|---|---|
| Initial (hard) capping | LPAR | LPAR share of CPC capacity | Any | Yes | No | SE/HMC |
| LPAR Absolute capping (zEC12 GA2 and later) | LPAR | Fractional #processors | | Yes | No | |
| LPAR Group Absolute Capping (z13 GA2 and later) | Group of LPARs | Fractional #processors | | Yes | No | |
| Defined capacity (DC, soft capping) | LPAR | MSU (4HRA) | CP | No | Yes | |
| LPAR group capacity (GC, soft capping) | Group of LPARs | MSU (4HRA) | | No | Yes | |
| Absolute MSU Capping | LPAR or Group | MSU | | Yes –Only for CPs– | No | SE/HMC + IEAOPT |
| Resource group or Tenant Resource Group (*) capping | Groups of service classes or Tenant Report Classes in Sysplex or per LPAR | Unweighted CPU SU/sec, fraction of LPAR share, fractional #CPs, or "accounted MSU" | CP or combined | N/A | No | WLM Service Definition |
| Logical configuration | LPAR | Integer #processors | Any | Yes | (Yes) | HMC+ OS |

# Which capping techniques may be combined?

−See next chart for legend−

| Type of capping ➜ | Initial (hard capping) | LPAR Absolute capping | LPAR Absolute group capping [2] | Defined capacity [1] | LPAR group capacity [1,2] | Resource group capping | Tenant Resource Group capping (*) |
|---|---|---|---|---|---|---|---|
| Initial (hard capping) | | + | + | - [3] | - [3] | + | + |
| LPAR Absolute capping | | | + | + | + | + | + |
| LPAR Group Absolute capping [2] | | | | + | +[2] | + | + |
| Defined capacity [1] | | | | | + | + | + |
| LPAR group capacity [1,2] | | | | | | + | + |
| Resource group capping | | | | | | | + [4] |
| Tenant Resource Group capping(*) | | | | | | | |

IBM Z

# Legend for *Which capping techniques may be combined?*

1) Includes ABSMSUCAPPING=NO and ABSMSUCAPPING=YES

2) Any LPAR can be defined to one group at most:
   Therefore, the group used for LPAR absolute group is the same as the group capping group

3) When initial capping is in effect, WLM cannot control capping:
   - Any defined capacity limit, if specified, will be ignored
   - The LPAR will not join a capacity group, or leave it, respectively.

4) Resource group and Tenant Resource Group capping may be combined within a service definition, but
   - A service class that is associated with a resource group cannot be assigned a Tenant Report Class. In other words:
     Any work unit may be capped through a resource group, or tenant resource group, but never both.

# More granular resource controls

- ❑ Purpose: provide more granular control over CPU and memory consumption by workload

- ❑ Initial Focus: High demanding workloads that run on speciality engines like Java batch, SPARK, analytics, and zCloud workloads

- ❑ New controls:
  - ➤ Honor Priority by service class
  - ➤ Memory Limit for resource groups

| Function \ z/OS Release | z/OS V2.2 | z/OS V2.1 |
|---|---|---|
| WLM/SRM Support | OA50845 | OA50845 |
| RMF Reporting Enhancements | OA50760 | OA50760 |
| z/OS Supervisor Support | OA50953 | OA50953 |
| z/OS RSM (Real Storage Manager) | OA51171 | OA51171 |
| z/OSMF Reporting Enhancements | PI71118 | PI71084 |

# New Resource Controls: Honor Priority

```
--------------------------------------------------------------------
                      Modify a Service Class           Row 1 to 4 of 4
Command ===>  _____

Service Class Name . . . . . : FAST
Description  . . . . . . . . . Velocity=80 goal
Workload Name  . . . . . . . . STCWORK      (name or ?)
Base Resource Group  . . . . .              (name or ?)
Cpu Critical . . . . . . . . . NO           (YES or NO)
I/O Priority Group . . . . . . NORMAL       (NORMAL or HIGH)
Honor Priority . . . . . . . . NO           (DEFAULT or NO)

Specify BASE GOAL information. Action Codes: I=Insert new period,
E=Edit period, D=Delete period.

         -- Period --  ---------------- Goal -----------------
Action   #  Duration   Imp.  Description
  __     1  2000000     2    Average response time of 00:00:01.000
  __     2  2000000     3    Average response time of 00:00:10.000
  __     3             4    Execution velocity of 80
****************************** Bottom of data ***********************
```

- ❑ Specifies whether work in this service class is exempted from default IFAHONORPRIORITY and IIPHONORPRIORITY processing
  - ➢ Also for Service Class Overrides
  - ➢ Limitation to specialty engines enforced collaboratively with zIIPs

- ❑ Usage
  - ➢ Some zIIP work may be very latency sensitive and require to be dispatched quickly .
    - ❍ Namely some DB2 work, such as prefetch SRBs.
    - ❍ zIIP capacity may be constrained but CP capacity might be available to help
  - ➢ Recommendation:
    - ❍ At the system level (IEAOPTxx) specify or default to IIPHONORPRIORITY=Yes to allow CPs to help zIIP work.
    - ❍ Use the service class specific HonorPriority=No to selectively exclude work from receiving help.
      - ❑ Examples could be SPARK or Java batch that you do not want to be processed on general purpose processors

* Statements regarding IBM future direction and intent are subject to change or withdrawal, and represent goals and objectives only.

IBM Z                                    © 2017 IBM Corporation                                    28

# New Resource Controls: Memory Limit for Resource Groups

```
-----------------------------------------------------------------------------
                          Modify  a  Resource  Group
Command  ===>  _____

Enter  or  change  the  following  information:

Resource  Group  Name  .  .  .  .  . :  MEMLIMIT
Description  .  .  .  .  .  .  .  .  .  Provide  a  memory  limit  for  work

Define  Capacity:
__      1.   In  Service  Units  (Sysplex  Scope)
        2.   As  Percentage  of  the  LPAR  share  (System  Scope)
        3.   As  a  Number  of  CPs  times  100  (System  Scope)
Minimum  Capacity  .  .  .  .  .  .  .  _____
Maximum  Capacity  .  .  .  .  .  .  .  _____

Memory  Limit  (System  Scope)        24_____       GB
```

❑ Specifies the maximum amount of memory that address spaces associated with the resource group through classification may consume on the local system (System Scope)

  ➢ The attribute is specified as absolute value in GB in the range 1 – 99,999,999.

  ➢ Also for Resource Group Overrides

❑ Memory limit enforced collaboratively with RSM

# How Memory Pools work

❑ A Resource Group definition, as part of the WLM service definition, has Sysplex scope.

❑ All address spaces get classified by WLM and get a Service Class assigned
  ➢ If the assigned Service Class is associated to a resource group and the resource group specifies a memory limit, WLM notifies RSM to connect the address space to a memory pool

❑ When the first address space connects to a memory pool, RSM creates that pool. The pool name is equal to the resource group name.

❑ An RSM memory pool represents only a logical pool (=upper limit).
  ➢ It is not dedicating or reserving real storage.
  ➢ A memory pool can specify a size exceeding the real storage of a given system

❑ All work in the system is managed towards the global pool (total real storage).
  ➢ Address spaces connected to a memory pool are also subject to the pool limit.
    ○ When the pool limit is approached self-stealing is initiated to keep the number of frames within the limit.

❑ Address spaces may be temporarily deferred, if  the pool limit would be exceeded by adding the space.

* Statements regarding IBM future direction and intent are subject to change or withdrawal, and represent goals and objectives only.

# Memory Pools: Current Limitations

❑ When the memory pool support becomes available the following limitations apply:

➢ An active address space cannot be reclassified to another defined memory pool.
  ○ The address space has to terminate in order to be reclassified, or it must be reclassified to the global pool and then reclassified again to the new pool.
  ○ Exception: initiator address spaces for a new job

➢ Resource Group memory pool limits cannot be decreased while it is defined into a policy.
  ○ The only way to decrease the limit is to activate a policy that does not have the resource group defined and then activate a policy that defines the pool with a smaller limit.
  ○ However, a pool can be dynamically increased via a new policy activation.

➢ Memory related Sysevents (such as STGTEST) are not memory pool aware

❑ IBM recommends that you use memory pools when it is required to limit memory consumption for new workloads such as Apache SPARK that provided guidance on how to operate them in a memory pool.
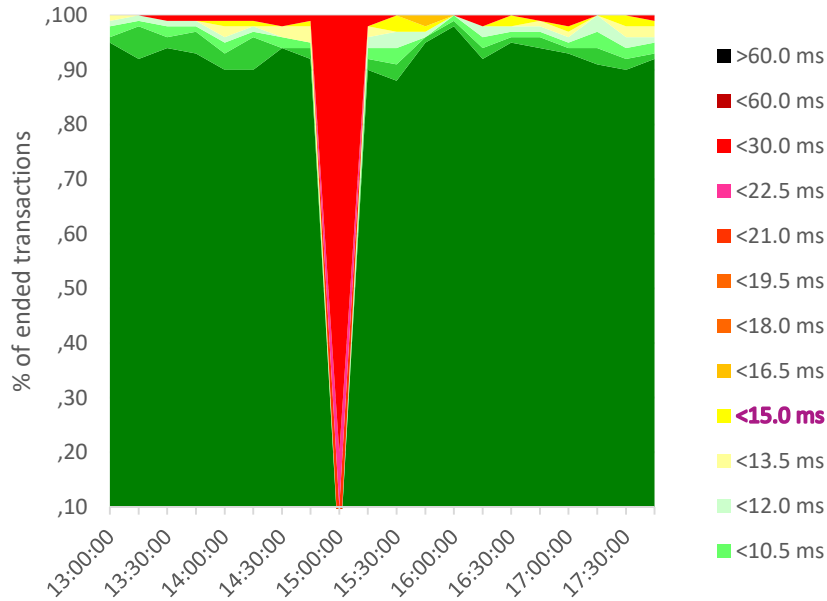
* Statements regarding IBM future direction and intent are subject to change or withdrawal, and represent goals and objectives only.

# WLM z14 Support

❏ Hiperdispatch Work Balancing Algorithm
  ➢ Takes memory affinity for high storage consumer into account
  ➢ If address spaces need to get splitted up the topology is considered and TCBs are moved to adjacent affinity nodes

❏ z/OS Freemained Frame management is extended to high frames (sometimes referred to as authorized storage)

| z/OS release / Function | V2.3 | V2.2 | V2.1 |
|---|---|---|---|
| WLM/SRM Support for IBM z14 | + | OA50144 | OA50144 |
| SRM Support for RSM Performance | + | OA51904 | OA51904 |

\* Statements regarding IBM future direction and intent are subject to change or withdrawal, and represent goals and objectives only.

# z/OS 2.3 Preview

❑ Use XML Format to save service definition!!

❑ Shorter Response Time Goals

➤ The current lower bound of 15 milliseconds for a response time goal is replaced by one millisecond allowing to specify meaningful goal values for very fast running transactions.

❑ Routing Enhancements for Softcapping

\* Statements regarding IBM future direction and intent are subject to change or withdrawal, and represent goals and objectives only.

# z/OS 2.3: Use XML Format to save Service Definition

❑ For several releases WLM has supported to store a service definitions in **XML format**
  ➢ z/OSMF WLM task
  ➢ ISPF Administrative Application: "Save as XML"…

❑ XML format avoids particular problems with the ISPF tables format, namely coexistence behavior, when a new functionality level needs to be introduced, and the number of table columns needs to be extended.
  ➢ **For example, OA47042  introduces such a change**
  ➢ **And there is more to come!**

❑ **Recommendation:**
  **Use the XML-format for your WLM service definition data sets**

❑ As of z/OS 2.3 XML-format is the default for the "Save" and "Save as" actions in the ISPF Administrative Application*

# z/OS 2.3: Shorter Response Time Goals

Illustrates an exemplary response time distribution for a service class with a 15 msec response time goal and very fast running transactions. During the time period, approximately 95% of the transactions complete with an average response time below 7.5 msec (dark green area). Thus, the service class considerably overachieves its goal allowing it to run at a low dispatching priority. This makes the service class susceptible to workload spikes where the average response time of all transactions rises far above the goal (red area with the outlier around 15:00).

A more meaningful goal of 7 milliseconds will cause a better balanced response time distribution as illustrated in the following figure.
The workload always runs at reasonable dispatch priority preventing bad response times even with workload spikes.

IBM Z

# z/OS 2.3: Shorter Response Time Goals

❑ The WLM Administrative Application is enhanced to allow the definition of service period response time goals below 15 milliseconds which is the minimum response time goal for z/OS V2R2 and before.
  ➢ The new minimum goal value is 0.001 seconds (one millisecond) and can be defined for base goals or when overriding attributes for a service class.
  ➢ When specifying an average response time goal, the total response time can be between 0.001 seconds (one millisecond) and 24 hours.

❑ Change of Functionality Level:
  ➢ With z/OS V2R3, the WLM Application level displayed on the Definition Menu panel changes to LEVEL035. As soon as at least one service class period is defined with a response time goal below 0.015 seconds (15 milliseconds), the Functionality level is raised to LEVEL035.
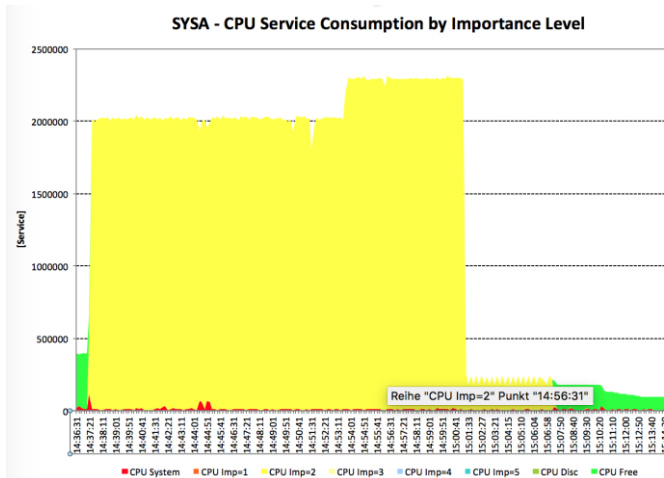
```
IWMAP34          Average response time goal

Enter a response time of up to 24 hours for period 2

Hours . . . . . 00  (0-24)
Minutes . . . . 00  (0-99)
Seconds . . . . 00.005  (0-9999)

Importance . . . 2  (1=highest, 5=lowest)
Duration . . . 10000      (1-999,999,999, or
                           none for last period)


F1=Help     F2=Split     F5=KeysHelp  F9=Swap    F12=Cancel
```

```
IWMAP35          Response time with percentile goal

Enter a percentile and response time goal for period 3

Percentile . . 90  (1-99)

Hours . . . . . 00  (0-24)
Minutes . . . . 00  (0-99)
Seconds . . . . 00.008  (0-9999)

Importance . . 2  (1=highest, 5=lowest)
Duration . . . 10000      (1-999,999,999, or
                           none for last period)


F1=Help     F2=Split     F5=KeysHelp  F9=Swap    F12=Cancel
```

\* Statements regarding IBM future direction and intent are subject to change or withdrawal, and represent goals and objectives only.
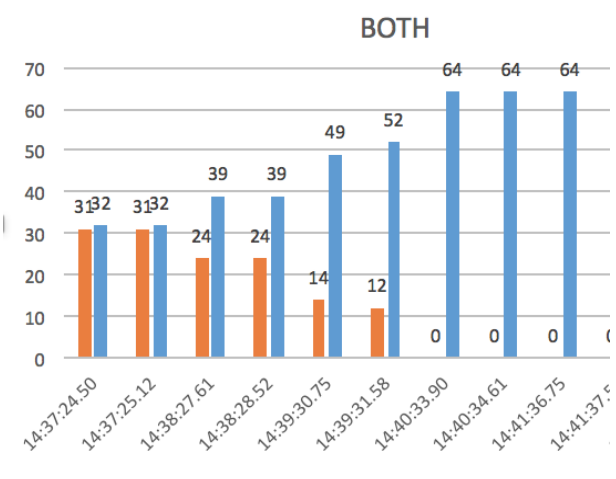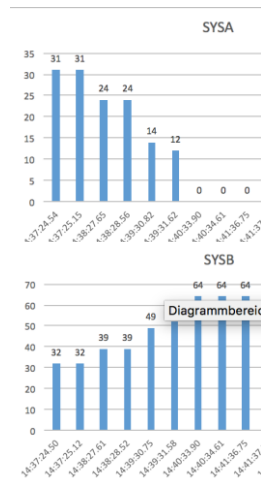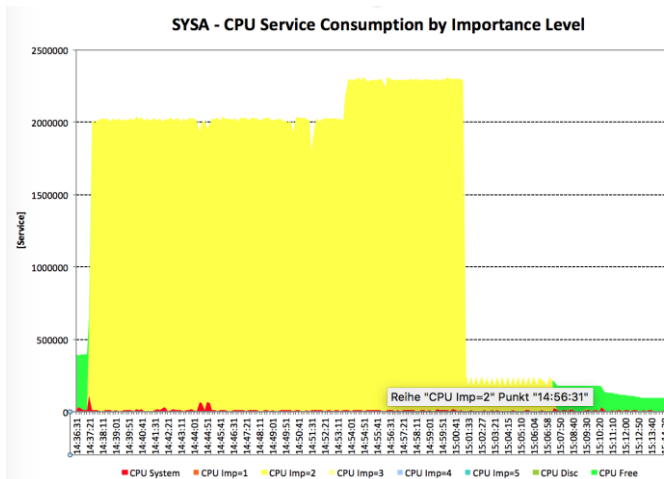
# z/OS 2.3: Routing Enhancements for Softcapping

❑ WLM Sysplex Routing services like IWMSRSRS and IWM4SRSC base their recommendation on the free and displaceable capacity of the systems in the Sysplex (3 minutes rolling average of actual capacities).

➢ This might result in routing work to a system which will be capped shortly thereafter due to Defined Capacity Limit or Group Capacity Limit.

❑ WLM is enhanced to take the capping limits into account when the free/displaceable capacity is determined. WLM will calculate the estimated time to capping for a system. The closer the system is to capping the more the available capacity will be reduced by the specified limit and influence the routing recommendations to send less work to the system.

# z/OS 2.3: Routing Enhancements for Softcapping

| IEAOPTxx RTCAPLEADTIME=nn | |
|---|---|
| <u>0</u> | The default behavior is as today: capping of this system will not be considered in advance. |
| [1-60] | Specifies the time in minutes, how long in advance an upcoming soft capping should influence WLM's Sysplex routing recommendations. When the estimated time to capping is less than n minutes WLM will consider the upcoming soft capping in its routing recommendations. |