# Oracle® Retail Planning Application Server (RPAS) on IBM AIX®

Version:

# RPAS Architecture

**IBM**®

Oracle Retail
9.x – 16.x

Classic Client    NOT: AP 14.1 or later

Select only one option

Fusion Client

Oracle Retail
13.0.4 and later

RPAS Classic Client

Client

WAN option: Terminal Server

ABC Intranet

RPAS Classic Client

Client

*Example Only!*

RPAS 16.0

WEB server

GBit Ethernet

App 1 RPAS Server

App 2 RPAS Server

App n RPAS Server

SAN

Domains

Domains

Domains

Copyright IBM, 2019

Client

WAN option: Terminal Server

ABC Intranet

WebLogic App. Server (Fusion Client)

WebLogic Server 12c (12.2.1)

IBM XLC 12.1 compiler used in 16.0
AIX 7.1 TL3 SP1 or AIX 7.2 TL1 SP1 and later.

**Note** that RPAS Server and WebLogic App. Server can be deployed on different OS'.

Client

GBit Ethernet

App 1 RPAS Server

App 2 RPAS Server

App n RPAS Server

SAN

Domains

Domains

Domains

# RPAS IO Characteristics - Observations

- RPAS IO is very random in nature with high write content

- Typically small average read IO sizes in the range of 30K to 80K; write IO sizes are typically larger, but are significantly influenced by AIX file cache size in comparison to active data.

- Low latency storage is essential – block storage & fibre connectivity

- Some sequential behavior, for example during "optimize domain" or workbook open, "auto-save" or when a copy of workbook data is made

- IO rates > 200,000 IOPS and IO throughput > 16GB/s during batch processing were observed in large environment – limited by SAN storage

- RPAS domains / workbooks are each made up of many files which need to be modified or re-created during the batch cycle. This behavior can lead to significant contention on JFS2 file system meta data control structures.

- RPAS applications can create a significant number of temporary files during processing → RPAS temporary files should not go to default "/tmp". Set the variable "TMPDIR" in the RPAS user environment to point to a fast SAN storage location. Regular maintenance of that directory is required as well.

# RPAS Data Layout for Optimal I/O Performance – SAN

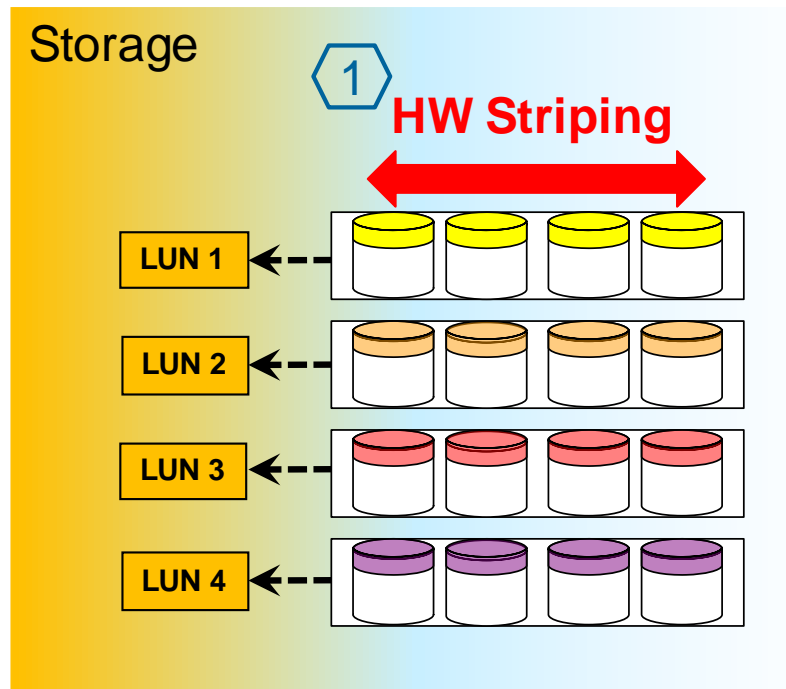**Example (with spinning disks):**

1. Use RAID-10 to create striped LUNs (in AIX seen as hdisks)
   - 1 LUN per RAID array per application
   - Each LUN is spread across 4 drives

### RAID-5 vs. RAID-10 Performance Comparison

| I/O Profile | RAID-5 | RAID-10 |
|---|---|---|
| Sequential Read | Excellent | Excellent |
| Sequential Write | Excellent (*1) | Good |
| Random Read | Excellent | Excellent |
| Random Write | Fair | Excellent |

RPAS IO characteristic
→ RAID-5 not a good option if using spinning disk



Storage

1 **HW Striping**
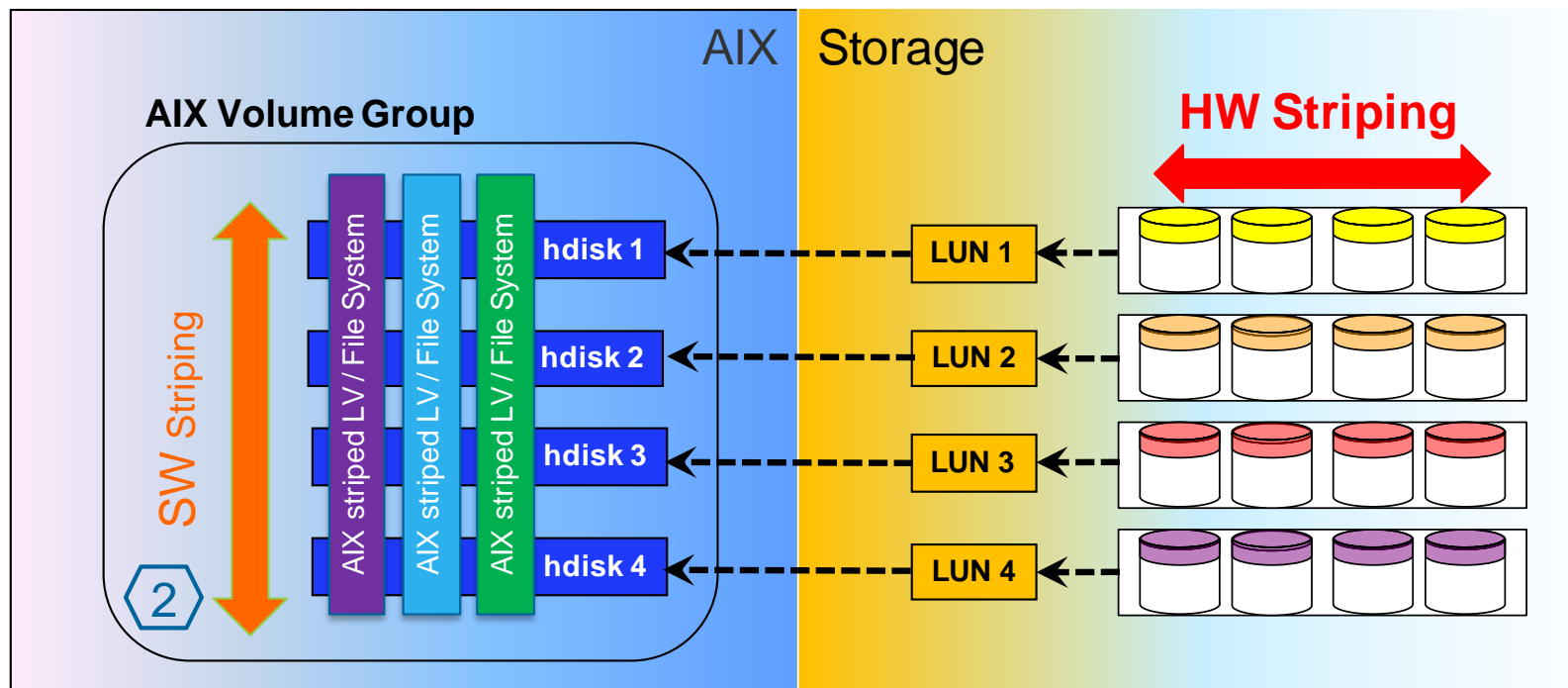
LUN 1
LUN 2
LUN 3
LUN 4

*1 – Assumes optimizing SAN storage sub-system!

# RPAS Data Layout for Optimal I/O Performance - AIX
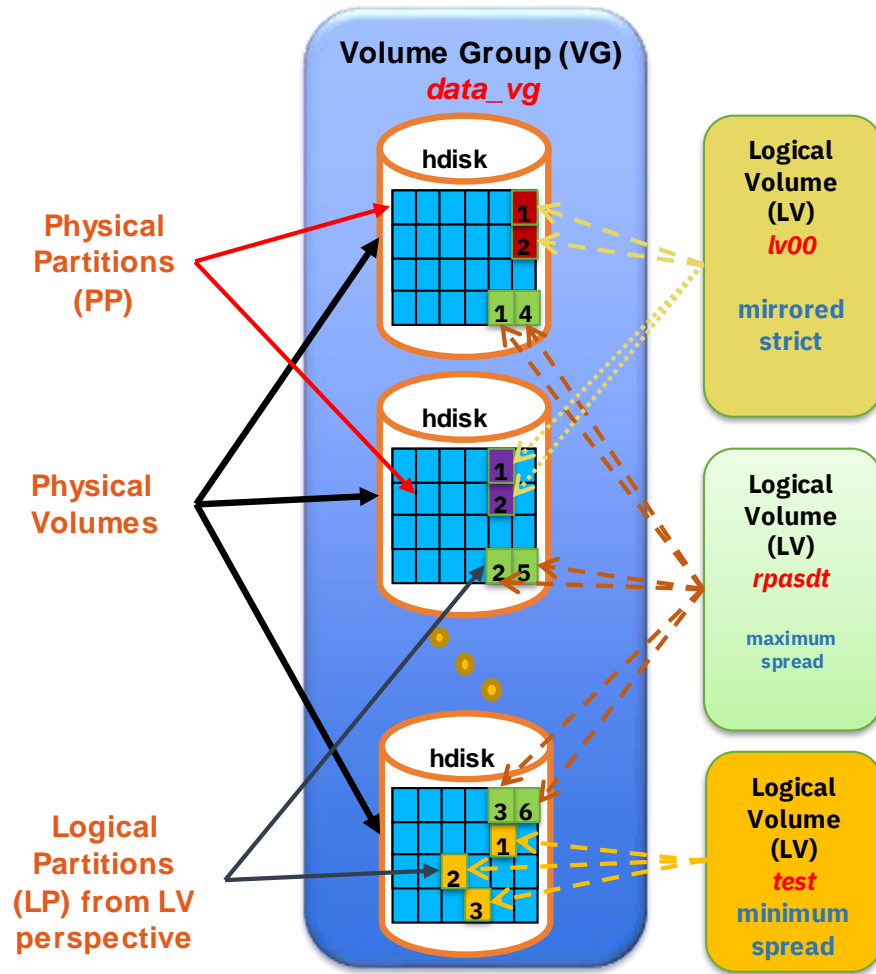
**Example**

(2.) Stripe or spread individual objects across multiple LUNs (hdisks) for maximum distribution
  – Each object is spread across 4 LUNs, each from different array (16 drives)



For IBM FlashSystem or A9000R utilize 8 or 16 LUNs of equal size per AIX VG.

**Note:** IBM Spectrum Scale (GPFS) for file systems hosting RPAS data can be very challenging to tune and is typically not recommended.

# Review - AIX Storage Logical Volume Management



- One hdisk belongs to 0 or 1 *VG*

- The size of an hdisk, SAN LUN for example, can be increased via SAN feature as needed and the *VG* increases in size accordingly

- An AIX system has 1 or more *VG*

- A *VG* contains *LP* from 0 or more *LV*

- A *LV* only contains *PP* from a single *VG*; it can not spread over multiple *VG*

- A *LV* can be increased in size up to the number of *PP* in the *VG*; the size of a *LV* can also be decreased, but only unused *LP* can be freed up

- *LV* can be created without software RAID, with RAID-0 and RAID 0+1

- *LV* can be used "raw" or a JFS, JFS2 file system is created on top of it

# Software Striping with AIX – LV striping & PP spreading

## Stripe using Logical Volume (LV)

- Create Logical Volume with the striping option : **mklv –a e –S <strip-size> ...**

- Strip size should be 1MB

- In VSCSI environments strip size should be 256K (VSCSI not recommended for RPAS)

- In NPIV / direct attached SAN storage environments a strip size of 1MB supports maximum throughput for sequential IO

Be aware of the implications LV striping has on the options to grow LV / FS size. SAN feature to grow LUN sizes dynamically is important. Key benefit is that files larger than 1MB will be striped over multiple hdisk.


## PP striping (AKA spreading)

- Create a Volume Group with a small PP size – <= 8MB – but keep in mind that scalable VGs with > 60k PP are significantly slower to create and maintain. (16MB PP size is alternative)

- Choose a **Scalable Volume Group : # *mkvg –S –s <PPsize> ...***

- Create LV with "***Maximum range of physical volume***" option to spread PPs on different hdisk in a round robin fashion :  # **mklv –a e –e x ...**

Challenge with PP striping is, that some "hot" data files in a domain / workbook can be smaller than a full PP stripe, so all IO to that file would go against only a subset of hdisks. PP striping increases risk of unbalanced IO if VG expansion is done by adding new LUNs. (reorgvg required !)

# LV striping with AIX – additional comments

- **The LV is stripped over *N* hdisks (typically all hdisks in a VG) from one VG.**

- **The LV space allocation can only grow in multiples of *N* times the *PP* size. Example: *N*=8, *PP* size=1GB → 8G, 16G, 24G, …**

- **A file system (FS) on top of a stripped LV should always grow with the same increments (*N* \* *PP*) so no space is wasted.**

- **If <span style="color:red">any</span> of the *N* hdisks runs out of free PP, attempts to grow the LV will fail. There are two options to resolve this:**
  1. Grow the underlying LUN(s) dynamically via SAN methods and discover the new size via "chvg –g <VG name>". Minimum size increase for each LUN is *PP* size. *Preferred Option!*
  2. Add another *N* hdisks (LUNs) to the VG and expand the LV to those new LUNs. This is called a "stripe column". The SA manually adds the hdisks to the VG and then expands the LV accordingly.

- **Option 1. is preferred as it requires no changes in SAN configuration and host mapping. It needs to be verified with the storage vendor if there are dynamic resize limitations. One such limitation is that a LUN in V7000, being part of a flashcopy relationship, can not be dynamically increased in size.**

- **Option 2. requires significant SAN changes and additional work by the AIX admin. Adding LUNs also impacts for example flashcopy configurations.**

# PP spreading with AIX – additional comments

- *PP* size needs to be carefully planned as LV/FS size is limited to 32k PP and management of VG with more than 60k PP becomes slow.

- The LV is PP spread over *M* hdisks (typically all hdisks in VG) from one VG.

- The LV space allocation can grow in multiples of *one* PP size, but should grow in multiples of *M* PP sizes for balanced IO distribution.

- A file system (FS) on top of a PP spread LV should always grow with the same increments (1+ or *M* * *PP*) so no space is wasted.

- If any of the *M* hdisks runs out of free PP, AIX skips that hdisk in the round-robin allocation from eligible hdisks in the VG. If no eligible hdisk has a free PP, further attempts to grow the LV will fail.

- There are two options to resolve this:
    1. Grow the underlying LUN(s) dynamically via SAN methods and discover the new size via "chvg –g <VG name>". Minimum size increase per LUN is *PP* size.
    2. Add another *K* hdisks (LUNs) to the VG and then execute "reorgvg <VG name>" to redistribute the allocated space evenly over all hdisks in the VG. *K >= 1.*

- Option 2. requires significant SAN changes and additional work by the AIX admin. Adding LUNs also impacts for example flashcopy configurations. "reorgvg" is an IO intensive operation.

# AIX JFS2 File Systems for RPAS

**IBM®**

- Pre-create the LV with the correct characteristics and then create the JFS2 file system on top of that LV with INLINE log. This may require the use of CLI.

- Limit the number of domains in one file system (mount point) to between 8 and 10.

- Create multiple JFS2 file systems within the same VG and then distribute domains by decreasing size round robin over the file systems. Master domain should be in a dedicated file system.

- Separating workbooks from their domains into their own FS', as well as dedicated VG(s), can be beneficial. Separation allows optimization for flashcopy and remote copy relationships as you can manage the VGs differently; for example, protect only master/local domain data, but not workbooks.

- Utilize "INLINE" JFS2 logs for JFS2 file systems created on top of striped LVs

- For JFS2 file systems on top of LV with "PP spreading" you can utilize "INLINE" JFS2 logs as well. Alternatively, a dedicated JFS2 redo LV per JFS2 file system with 256KB or 1MB strip size, 1 PP allocated from each hdisk in the VG, can provide potentially better performance and better JFS2 log IO distribution.

- RPAS file systems should be mounted with the "noatime" mount option; do NOT utilize "DIO" or "CIO". The latter will corrupt your data!

If multiple JFS2 file systems are needed, utilize the "copyDomain" or "moveDomain" utilities to distribute RPAS domains round robin over the file systems starting with the largest / most complex domain. **Do not utilize "symbolic file system links" as of RPAS version 15.**

# AIX JFS2 File Systems for RPAS – an example with striped LV

- **Change device settings; queue_depth needs to be verified with storage vendor and FC adapter settings need to be updated accordingly, especially num_cmd_elems.**
  - chdev -l hdisk35 -a max_transfer=0x100000 -a queue_depth=256
  - ...

- **Create the scalable volume group with 1GB PP size. With a goal of less than 60k PP in a VG this would support a VG size up to 50TB and a maximum FS size of 32TB.**
  - mkvg -S -s 1024 -y rpasvg_1 hdisk35 hdisk36 hdisk37 hdisk38 hdisk39 hdisk40 hdisk41 hdisk42

- **Create two striped LV with stripe width of 8, 1MB strip size, maximum number of 6000 LP in previously created VG. Note the randomization in hdisks to improve concurrent sequential IO.**
  - mklv -a e -c 1 -C 8 -S 1M -t jfs2 -x 6000 -y rpasd01lv  rpasvg_1 3000 hdisk35 hdisk36 hdisk37 hdisk38 hdisk39 hdisk40 hdisk41 hdisk42
  - mklv -a e -c 1 -C 8 -S 1M -t jfs2 -x 6000 -y rpasd02lv  rpasvg_1 3000 hdisk39 hdisk42 hdisk38 hdisk41 hdisk36 hdisk40 hdisk37 hdisk35

- **Create JFS2 file systems with "INLINE" jfs2 log and "noatime" mount option on previously generated LVs using all available space in each LV.**
  - crfs -v jfs2 -p rw -a logname=INLINE -a options=noatime -A yes -d rpasd01lv -m /rpasd01
  - crfs -v jfs2 -p rw -a logname=INLINE -a options=noatime -A yes -d rpasd02lv -m /rpasd02

# RPAS on AIX – Memory

- A **large AIX file system cache** is typically required to reduce the IO workload generated by RPAS against the SAN. Typical memory configurations for a LPAR running RPAS are between 12GB and 16GB per processor.

- **Do not over configure the maximum memory** setting for the LPAR definition, especially on multi-CEC environments, to achieve best memory affinity between CPU and memory.

- **Utilize AIX 64k memory page size** for RPAS - but closely monitor AIX "psmd" daemon and page conversions between 4k <-> 64k pages. You want to ensure that you have sufficient 64k pages available to support RPAS binaries and their stack / local data. The use of 64KB pages is set in the RPAS user environment via:

  ```
  export LDR_CNTRL=DATAPSIZE=64K@TEXTPSIZE=64K@STACKPSIZE=64K
  ```

- Set "**export MALLOCOPTIONS=pool:0x30000000,buckets,no_mallinfo**" in the RPAS user environment to reduce amount of CPU utilized by RPAS.

**Ensure AIX 7.2 fix for APAR IJ09762: "MALLOCOPTIONS=POOL,BUCKETS,NO_MALLINFO CAUSES MEMORY BLOAT" is applied!**

**Tip:**
nmon + "M" reports memory by page size

# RPAS on AIX – Memory – "lrud"

- The AIX lrud daemon is responsible to free up physical memory pages to support new memory requests for the application data or for IO to disk.

- The lrud can use a significant amount of CPU resources and under large memory pressure become a bottleneck. In the context of RPAS, cached IO to disk can significantly negatively impact especially batch throughput.

- To reduce the impact of lrud during batch workloads you want to utilize specific JFS2 mount options during different batch phases. Those mount options can be dynamically changed without application outage.

- Recommended mount options to test with – "noatime" for all FS:

  - **Interactive workload**, use default mount options

  - **Batch workload which reads domain data in and writes modified domain data** out, without accessing the same data again, use "rbr,rbw". File system with master domain data should be always mounted with default flags (no rbr or rbw)!

  - **Batch workbook creation / update** – file systems for master/local domains are mounted with default options. File systems with workbooks are mounted with "rbr,rbw".

**Tip:**
*mount –o remount,...* allows to dynamically change the mount options for a file system.

# RPAS on AIX – CPU

- Utilize the **highest frequency POWER CPU** to best support the single threaded characteristic of RPAS applications.

- **Utilize Shared Multi Threading (SMT)** for RPAS workloads for best throughput; if you do not have sufficient independent domains to utilize all logical processors, then reducing SMT mode accordingly can be beneficial. Do not reduce below SMT=2 on POWER9 processor.

- Set the LPAR **CPU entitlement to a value reflecting typical CPU usage**, at minimum to what is regularly required to support interactive user workload

- **Do not over configure the number of virtual processors**, especially on multi-CEC environments, to achieve best CPU to memory affinity and optimize processor throughput

- It is typically beneficial to configure virtual processor (VP) folding to keep one or two additional VP active above what AIX calculates. This change can help in addressing spikey workload demand better. ( *schedo –o vpm_xvcpus=1(2)* )

- The configured RPAS parallel processing limits should be set so that the average run queue, as reported by vmstat, is close to the number of logical processors in the LPAR. Typically this leads to a ratio of 1.1 to 1.5 of parallel RPAS processes configured to number of logical processors. A higher ratio required to keep all CPU busy in the LPAR typically indicates a more limiting IO subsystem.

# RPAS on AIX – IO Tuning Options

- To reduce the risk of running out of JFS2 IO buffers the following should be set:
**ioo –p –o j2_dynamicBufferPreallocation=256**

- To reduce the impact of the syncd process, the following settings can be implemented; this will result in syncd be more "lazy", which is typically not a reason for concern with the RPAS application as an unexpected outage typically requires a full restore of the data anyhow.
**ioo –p –o j2_syncPageCount=512**
**ioo –p –o j2_syncPageLimit=4096**

- During RPAS batch processing AIX is somewhat overly aggressive in writing dirty pages back to disk. RPAS has the tendency to re-write a data block several times during processing and the "early" write of AIX can result in multiple writes of the same block to disk. Setting the following parameter can significantly reduce that "re-write" behavior. It is not recommended to set this parameter to "0" (disable) if not all file data can be kept in the JFS2 file cache.

**ioo –p –o j2_nPagesPerWriteBehindCluster=480**

# RPAS on AIX – Tuning Options (Advanced)

- RPAS parallel processing, and how files are handled, can lead to significant file fragmentation with default AIX settings. AIX internal improvements are under way, but until they are released, it is possible to influence AIX to pre-allocate more continues space for a file than the default. At this time it requires the use of kdb, should only be done by a skilled AIX admin and tested in non-production first (No WARRANTY!). Here the commands to execute in an AIX 7.2 environment – needs to be re-done after every reboot:

  ```
  #!/bin/ksh

  kdb -script <<EOM
  dw j2_tunables+90
  mw j2_tunables+90
  1e0
  .
  q
  EOM
  ```

- To reset to default value re-run commands above and replace "1e0" with "20".

- AIX 7.2 also provides the **defragfs** command to dynamically defragment file systems. Before using that command ensure that you are running AIX 7.2 TL3 release level.
(https://www.ibm.com/support/knowledgecenter/en/ssw_aix_72/com.ibm.aix.cmds2/defragfs.htm)

# RPAS on AIX – Miscellaneous

- **Use RAID-10** for production environments with spinning disks to support the highly random IO behavior of this workload; especially the large number of random writes. For SSD / Flash only environments RAID-5 is an option.

- **Spread IO** over all physical disks as evenly as possible. Applies to SSDs as well.

- Configure the RPAS application with **as many domains as possible**, while still meeting business requirements.

- Configure RPAS local domains so they are about **equal in size / (compute) complexity**

- **Keep workbooks small** for faster processing times

- RPAS has a limited/no transaction concept. This means, that **recovery** from a data error during batch processing or an unplanned outage, **typically requires a restore** from the last known good backup to bring the RPAS domains and workbooks back into a known good state.

- RPAS environments typically require a **SAN based "fast backup / restore" option** like "flash copy" in IBM storage solutions.

- Under the premise that RPAS has no transaction concept, you could disable JFS2 logging on all file systems with master / local domain data. This can significantly improve batch processing, but means that in case of an unplanned outage, or "bad data loaded", all related file systems need to be re-created and then the data restored. Those FS' should also NOT be automatically mounted at AIX boot time as "chkfs" will take a long time and likely fail.

# RPAS on AIX – Miscellaneous (continued)

It has been observed that unmounting of file systems with large amounts of data in the JFS2 file cache can take a significant amount of time. (part of reboot/shutdown as well).

- The latest AIX perfpmr packages (**ftp://ftp.software.ibm.com/aix/tools/perftools/perfpmr**) contain a small tool called "flushfile" written by Mathew Accapadi, IBM. "flushfile" run without parameters shows command options.

- The tool allows to "un-cache" content of any file from the AIX file cache and has proven to be much faster than the process used in unmount as it can be run in parallel on multiple directories under the same mount point. This could also be beneficial in preparation of a SAN based "flash copy" of VGs with RPAS domain data.

- In the context of RPAS, the following is an example approach to release all cached pages belonging to files in three file systems with RPAS domain directories / files – note that this can create a heavy write workload as all dirty pages have to be written to disk before a page is removed from file cache!

```
#!/bin/sh
for wrkFS in rpasFS1 rpasFS2 rpasFS3; do
        for ldom in $wrkFS/*; do
                ./flushfile –v 1 –s 10000 –r –d $ldom &
        done
done
echo "Waiting for flushfile to complete ... \c"
wait
```

# RPAS on AIX – RPAS Fusion UI

- The Fusion web-based UI runs within a WebLogic server instance.

- Optimize communications between Fusion Application & RPAS server by setting the following environment variables:
  - *FC_TCPNODELAY=TRUE*
  - *RPAS_TCPNODELAY=ON*     *(RPAS 13.3, 13.3.1 and 13.4 or later)*

- For optimal UI response times, tune WLS Java garbage collection for low-latency *(-Xgcpolicy:gencon)*

- Set minimum and maximum Java heap sizes *(-Xmx, -Xms)* identically to reduce heap expansion/contraction times.

- Sizing the heap correctly requires monitoring of the JVM in order to minimize the time spent in garbage collection.

- A 4 GB Java heap is a good starting point for the WLS instance.

- **Issue:** Old RPAS Fusion client not being able to connect anymore
  https://www.ibm.com/support/knowledgecenter/en/SSYKE2_7.0.0/com.ibm.java.security.component.70.doc/security-component/jsse2Docs/disablesslv3.html

> **Tip:**
> Use IBM Support Assistant – Garbage Collection and Memory visualizer to analyze heap usage.

# RPAS - Backup / Recovery

- **(#1)** Enables recovery from data error and unplanned system outage during batch processing

- **(#2)** Includes verified global/local domains with all pre-built / old workbooks and provides:
  - Restore point to recover to before next successful **(#1)** backup
  - Source for persistent backup (tape) required for long term data retention / disaster recovery

- **(#2)** typically should not overwrite **(#1)** to enable quick recovery in case a data issue is discovered later during end-user processing the next day.

- Depending on application, daily batch may be required and above approach could be applied to the daily batch cycle as well.

- The daily persistent backup typically would be a differential backup in contrast to the typical full backup after weekly batch processing.

# RPAS - Backup / Recovery (Continued)

- If the step "Backup All Domain Data" is implemented via SAN functionality like "FlashCopy", then it needs to be ensured that before taking that copy that there is no data left in the AIX file cache which has not yet been written to disk.

- There are 2 options to achieve that before taking the FlashCopy in the SAN:

  1. "unmount all related file systems"

  2. "freeze" (JFS2 only; *chfs –a freeze=<freeze time> <FS>*) all related file systems and wait for all dirty pages to be flushed to disk

- In the context of 1., you want to evaluate the use of the "flushfile" utility mentioned earlier in this presentation as it can significantly reduce the time to successfully un-mount a file system.

- The benefit of option 2. is that domain data stays in cache and is ready for batch processing after the freeze is removed. A potential challenge is that there is no indication if / when all metadata / file blocks have been flushed to disk – freeze de-stages meta-data and dirty pages, but that can take some time depending on dirty pages in cache and IO sub-system performance. You could potentially monitor via "nmon" or "iostat" for any write IO to underlying hdisk(s).

# Notices and Disclaimers

# Notices and Disclaimers continued

- Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products about this publication and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products. IBM does not warrant the quality of any third-party products, or the ability of any such third-party products to interoperate with IBM's products. **IBM expressly disclaims all warranties, expressed or implied, including but not limited to, the implied warranties of merchantability and fitness for a purpose.**

- The provision of the information contained herein is not intended to, and does not, grant any right or license under any IBM patents, copyrights, trademarks or other intellectual property right.

- IBM, AIX, Easy Tier, IBM FlashCore, IBM FlashSystem, IBM Power Systems, IBM PureSystems, IBM Spectrum, IBM Spectrum Accelerate, IBM Spectrum Archive, IBM Spectrum Control, IBM Spectrum Protect, IBM Spectrum Scale, IBM Spectrum Storage, IBM Spectrum Virtualize, IBM Watson, IBM z Systems, IBM z14, OpenPower, POWER, Power Systems, PowerHA, PowerLinux, PowerVM, Real-time Compression, Storwize, and XIV are trademarks of International Business Machines Corporation, registered in many jurisdictions worldwide.

- Other product and service names might be trademarks of IBM. A current list of IBM trademarks is available on the Web at "Copyright and trademark information" at www.ibm.com/legal/copytrade.shtml.

- The following terms are trademarks of other companies:
  - Linux is a trademark of Linus Torvalds in the United States, other countries, or both.
  - Intel is a trademark of Intel Corporation or its subsidiaries in the U.S. and/or other countries.
  - UNIX is a registered trademark of The Open Group in the United States and other countries.
  - Oracle and Java are registered trademarks of Oracle and/or its affiliates.
  - Other company, product, or service names may be trademarks or service marks of others.