In the past few releases, DB2 for z/OS has steadily increased its exploitation of zIIP speciality processors. In addition to providing significant potential for mainframe cost reduction, zIIP offload can have a positive (or negative) impact on DB2 query performance. This presentation will review DB2 zIIP usage by release, and provide the attendee with some practical guidance on zIIP capacity planning and usage monitoring.

IDUG

Leading the DB2 User
Community since 1988

**IDUG DB2 Tech Conference**
Prague, Czech Republic | November 2014

🐦 #IDUG

2

## Acknowledgements

- John Lemmon, Triton Consulting

- Adrian Burke, IBM
- Florence Dubois, IBM
- Willie Favero, IBM
- James Gill, Triton Consulting
- Martin Packer, IBM

3

## Topics

- Introduction
- What Are zIIPs and How Can They Help?
- How Does DB2 Exploit zIIPs?
- Performance Considerations
- zIIP Usage Measurement and Capacity Planning
- Case Study
- Summary

4

## Introduction

- DB2 consultant with Triton Consulting, based in the UK
- 26 years DB2 experience
  - Database Administration
  - Systems Programming
  - Application Development
- IBM Gold Consultant
- IBM Champion
- IBM White Papers, Redbooks, Flashbooks, etc
- IDUG Best Speaker and Past President

## Topics

- Introduction
- **What Are zIIPs and How Can They Help?**
- How Does DB2 Exploit zIIPs?
- Performance Considerations
- zIIP Usage Measurement and Capacity Planning
- Case Study
- Summary

# What Are zIIPs?

- zIIP – z/OS Integrated Information Processor
- IBM System z speciality processor (SP)
  - Able to execute work that has been redirected from a Central Processor (CP) under certain conditions
  - Work that runs on a speciality processor does not count towards z/OS Monthly Licence Charge (MLC) software fees
  - Other speciality processors include zAAP, IFL, ICF
- zIIP designed to support enclave SRB workloads

6

IDUG
Leading the DB2 User
Community since 1988

**IDUG DB2 Tech Conference**
Prague, Czech Republic | November 2014

🐦 #IDUG
7

## What Are zIIPs?

- zIIP is just a re-tasked CP in terms of physical processor unit (PU) chip
  - Always runs at full PU capacity, even if CEC is a "kneecapped" / "hobbled" sub-capacity model
  - SMF70NRM field provides info on CP/zIIP normalisation factor but DB2 SMF data already normalised
- Evolution
  - Introduced in 2006 with z9 – limit of 1 zIIP per CP
  - zEC12/zBC12 upped limit to maximum of 2 zIIPs per CP

Field SMF70NRM gives us a number which we can divide by 256 and see how much faster the specialty engines are than the CPs

IDUG
Leading the DB2 User
Community since 1988

IDUG DB2 Tech Conference
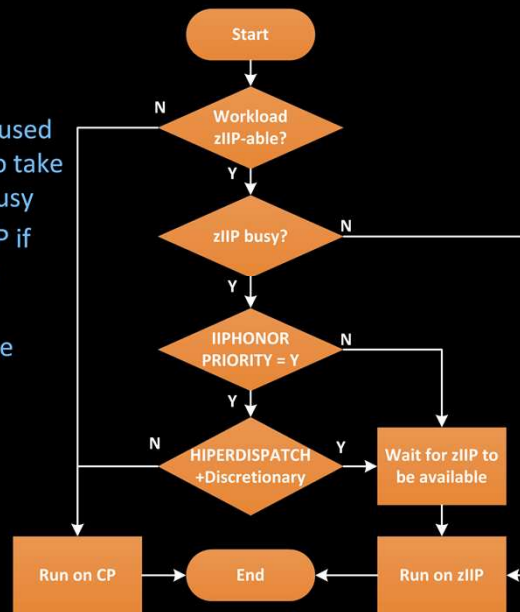Prague, Czech Republic | November 2014

#IDUG
8

## How Can zIIPs Help?

- Reduce Total Cost of Ownership (TCO)
  - One-off fee to purchase zIIP (costs much less than a CP), with ongoing potential savings in software MLC fees
  - Potential to further reduce costs for some footprint-licenced software (e.g. SAS) at next technology refresh (e.g. z196 to zEC12) if CP MSU total can be reduced
  - Provide additional capacity headroom, allow expensive mainframe CP upgrades to be avoided / deferred
- Performance?
  - Secondary potential benefit in some specific cases (e.g. CPU starvation)
  - Can degrade performance in some situations as well (see later)

Also allows next technology refresh (e.g. z196 => zEC12) to provide less total MSUs and therefore a reduction in some ISV software – e.g. SAS is footprint licence.

IDUG
Leading the DB2 User
Community since 1988

IDUG DB2 Tech Conference
Prague, Czech Republic | November 2014

🐦 #IDUG
9

## What Runs Where?

- IIPHONORPRIORITY PARMLIB option used to specify whether CPs are allowed to take on zIIP-eligible work if zIIPs are too busy
- YES – work can be directed back to CP if zIIP too busy (default and usually the recommended value)
- NO – work will wait for zIIP to become available (can manifest as NOT ACCOUNTED time)
- Beware combination of HIPERDISPATCH=YES and discretionary zIIP-eligible workload
  - Forces work to wait regardless of IIPHONORPRIORITY setting

Start

N ← Workload zIIP-able?
Y ↓

zIIP busy? → N
Y ↓

IIPHONOR PRIORITY = Y → N
Y ↓

N ← HIPERDISPATCH +Discretionary → Y → Wait for zIIP to be available

Run on CP → End ← Run on zIIP

Simplified – not all situations covered!

HiperDispatch is a workload dispatching feature found in the newest IBM mainframe models (the System z10 and IBM zEnterprise System processors) running recent releases of z/OS. HiperDispatch was introduced in February, 2008.

One of the engineering challenges with large SMP server designs is to maintain near-linear scalability as the number of CPUs increases. Performance and throughput do not double when doubling the number of processors. There are many overhead factors, including contention for cache and main memory access. These overhead factors become increasingly difficult to mitigate as the number of CPUs increases. The design goal for delivering maximum performance is to minimize those overhead factors. Each new mainframe model supports a higher maximum number of CPUs (up to 64 main processors in a single System z10 mainframe for example), so this engineering challenge becomes ever more important.

HiperDispatch helps address the problem through a combination of hardware features, z/OS dispatching, and the z/OS Workload Manager. In z/OS there may be tasks waiting for processing attention, such as transaction programs. Each task often requires access to memory. In a large SMP design such as System z, some CPUs are physically "closer" with faster access to cache memory that might hold supporting data for particular tasks. HiperDispatch exploits this fact and steers tasks to the CPUs most likely to have the fastest access to relevant data already in cache. If that particular CPU is busy, HiperDispatch will, at first, wait for it to

finish its other task, even if another less favorable CPU is idle. However, there are limitations to how patient HiperDispatch will be, as governed by Workload Manager goals. If z/OS Workload Manager senses that there's a risk the pending task will miss its service level (responding within a certain number of milliseconds to a user request for example), Workload Manager and HiperDispatch will send the task over to an idle CPU for processing, even if that CPU must fetch data from slower main memory.

IDUG
Leading the DB2 User
Community since 1988

**IDUG DB2 Tech Conference**
Prague, Czech Republic | November 2014

#IDUG
10

## What Runs Where?

- ZIIPAWMT PARMLIB option used to specify how often WLM checks to see if a zIIP needs help
  - Fixed at 3200 for HIPERDISPATCH=YES, so delay of up to 3.2 m/s if zIIP too busy and workload has to be directed back to a CP
- ZIIPMAXQL PARMLIB option used to specify maximum number of dispatchable units that will queue waiting for a zIIP processor
  - Default is 7
- IBM advice is not to change either of these unless advised by L2 support

**IDUG**
Leading the DB2 User
Community since 1988

**IDUG DB2 Tech Conference**
Prague, Czech Republic | November 2014

🐦 #IDUG

11

# Terminology

- zIIP-able
  - Workload that is eligible for zIIP offload

- zIIP-ed
  - zIIP-able workload that is actually executed on a zIIP

- Un-zIIP-ed
  - zIIP-able workload that ends up being executed on a CP

IDUG
Leading the DB2 User
Community since 1988

**IDUG DB2 Tech Conference**
Prague, Czech Republic | November 2014

🐦 #IDUG
12

**Topics**

- Introduction
- What Are zIIPs and How Can They Help?
- **How Does DB2 Exploit zIIPs?**
- Performance Considerations
- zIIP Usage Measurement and Capacity Planning
- Case Study
- Summary

IDUG
Leading the DB2 User
Community since 1988

IDUG DB2 Tech Conference
Prague, Czech Republic | November 2014

#IDUG
13

## DB2 zIIP Exploitation

- DB2 has steadily made more and more work zIIP-able from V8 onwards, and trend is expected to continue
  - Significant jump in exploitation in DB2 10 onwards, capacity planning becomes much more important!
- Some consistent themes for offloadable work
  - "New" workloads – distributed access, XML
  - Utilities
  - Async system tasks (prefetch, pseudo-deleted index clean-up)
- Good instrumentation is available
  - Both general and DB2-specific SMF records – see later

14

## DB2 zIIP Exploitation

- DB2 is currently the major (but not the only) zIIP exploiter
  - DFSORT
  - z/OS Communications Server (IPSec encryption, HiperSockets Multiple Write)
  - z/OS Global Mirror (aka XRC)
  - zAAP on zIIP is driving usage (WAS on z/OS, Java stored procedures)
  - Many ISV products / functions
- DB2 is the first exploiter to begin making performance critical tasks zIIP-able

Prefetch and Deferred Write processing : 100% (c. 70% of DBM1 SRB time)

IDUG
Leading the DB2 User
Community since 1988

**IDUG DB2 Tech Conference**
Prague, Czech Republic | November 2014

#IDUG
17

## Performance Considerations

- zIIPs are primarily intended as a means of reducing TCO by managing MLC costs, not a performance boost
- BUT… zIIPs can certainly have in impact on performance!
- Performance can sometimes be improved if you are driving your CPs very hard and workload is suffering from CPU starvation
  - Enabling offload to zIIP effectively increases overall compute capacity, thereby reducing CPU starvation and elapsed time
- Performance can often be degraded if you're driving your zIIPs too hard
  - Wait for zIIP, or wait to be redirected back to CP

17

**IDUG**
Leading the DB2 User
Community since 1988

**IDUG DB2 Tech Conference**
Prague, Czech Republic | November 2014

#IDUG
18

## zIIP Utilisation – How Much is Too Much?

- Many customers drive their CPs very hard (in the 90%-100% utilisation range), and that's OK
  - Entire infrastructure is designed to cope with very high utilisation levels
- zIIPs typically cannot be driven nearly as hard without significant performance issues
  - Aiming for average of 30-50% utilisation per RMF interval is good rule of thumb (peaks can of course be higher)
- DB2 10 made several performance-critical system tasks zIIP-able (and DB2 11 adds some more)
  - Performance impact of insufficient zIIP capacity is increasing with each release

IDUG
Leading the DB2 User
Community since 1988

IDUG DB2 Tech Conference
Prague, Czech Republic | November 2014

#IDUG
19

## zIIP Utilisation – How Much is Too Much?

- We take a triple hit if redirecting back to CP (un-zIIP-ed)
    - Delay caused by waiting for zIIP then going back to CP
    - Workload will run on CP, so contributes to 4HRA and could impact MLC fees
    - Work may take longer to execute if being redirected back to sub-capacity CP (zIIPs always run at full PU capacity)
- Disabling zIIP (i.e. taking offline from LPAR) can be an effective short-term measure
    - May have MLC implications!
    - Proper long-term solution is probably to purchase more zIIP capacity

# Specific Situations to Avoid

- DB2 10 and above, HIPERDISPATCH=YES, IIPHONORPRIORITY=YES, zIIPs too busy (averaging >> 30-50%)
  - Potential for significant slowdown in prefetch activity, leading to widespread DB2 performance issues
  - See Flo's blog: https://www.ibm.com/developerworks/community/blogs/22586cb0-8817-4d2c-ae74-0ddcc2a409bc/entry/december_17_2012_6_07_am3?lang=en
- zIIP-able work, HIPERDISPATCH=YES, IIPHONORPRIORITY=YES, performance critical workload incorrectly classified as discretionary, over-stressed zIIP
  - zIIP-able work will still be forced to run on zIIP, potentially leading to significant elapsed time issues
  - See Willie's blog: http://db2onlinehandbook.com/?p=566

## Specific Situations to Avoid

- Low-weight LPARs with low zIIP usage, sharing a single zIIP
  - Possible latency issues waiting for logical zIIP to be dispatched if used by other LPARs

## zIIP Capacity Planning is important

- DB2 10 and DB2 11 allow significantly more offloading
  - Time to get serious about zIIP Capacity Planning if you want to optimise Cost Reduction opportunities
- Understanding how IBM software pricing works isn't essential, but it will certainly help
- Tools
  - General SMF/RMF (30, 70, 72)
  - DB2 SMF (100, 101)
  - Other (SCRT)
- Make sure you set PARMLIB PROJECTCPU=YES, this directs z/OS to record the amount of work eligible for zIIP even if you don't have one installed

IDUG
Leading the DB2 User
Community since 1988

**IDUG DB2 Tech Conference**
Prague, Czech Republic | November 2014

#IDUG

24

## SMF 30 Data

- SMF 30 (Common Address Space Work / Activity)
  - 30.2 & 30.3 – Written at end of interval
  - 30.4 – Written at end of step
  - 30.5 – Written at end of job
- CPU times are in hundredths of seconds (.01 seconds)
- Useful zIIP fields
  - SMF30_TIME_ON_ZIIP – actual zIIP offload (zIIP-ed)
  - SMF30_TIME_ZIIP_ON_CP – zIIP-eligible workload redirected to CP (un-zIIP-ed)
- Broken down by LPAR, Address Space, WLM name, Service Class, Userid, Account-code…

# SMF 30 data can tell us...

- Current zIIP offload (zIIP-ed) broken down by Service Class
- SSYHI looks like it may contain DB2 in this example...

## SMF 30 data can tell us...

- To confirm, we can then chart the same data by PGM prefix
- Same data as shown on previous chart
  - DSN* and DXR* for DB2

## SMF 30 data can tell us…

- Overall zIIP utilisation
- Shown as
  - MSU (remember to normalise)
  - zIIP %busy (remember, ROT is to not exceed 35-50% busy)



28

IDUG

Leading the DB2 User
Community since 1988

**IDUG DB2 Tech Conference**
Prague, Czech Republic | November 2014

#IDUG

29

# SMF 30 data for DBM1 can tell us…

- If you're still on V9, measuring current zIIP utilisation won't give you much of a handle on what will happen in V10
  - Prefetch / deferred write not reported as IIPCP (un-zIIPed)
- Different approach needed
  - Up to 100% of DBM1 address space becomes eligible at v10+, so use SMF30 to measure total DBM1 MSU usage to get zIIP estimate



Job Interval MSUs by Product

IDUG
Leading the DB2 User
Community since 1988

IDUG DB2 Tech Conference
Prague, Czech Republic | November 2014
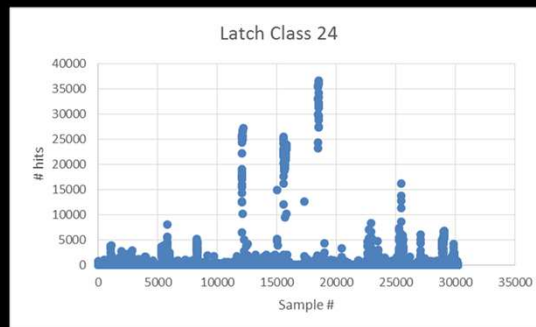
🐦 #IDUG
32

32

## DB2 Stats can tell us…

- DB2 statistics reports show zIIP offload from system address spaces
- DB2 10
  - Prefetch and deferred write processing for DBM1
  - PM30468 ensures zIIP-ed workload correctly reported under DBM1
- DB2 11
  - Log write, log prefetch MSTR
  - Index pseudo-delete clean-up daemons for DBM1

zIIP-ed offload

| CPU TIMES | TCB TIME | PREEMPT SRB | NONPREEMPT SRB | CP CPU TIME | PREEMPT IIP SRB |
|---|---|---|---|---|---|
| SYSTEM SERVICES ADDRESS SPACE | 0.120789 | 0.037248 | 0.027924 | 0.185961 | 0.000000 |
| DATABASE SERVICES ADDRESS SPACE | 0.403656 | 0.218449 | 0.009574 | 0.631680 | 0.000000 |
| IRLM | 0.000043 | 0.000000 | 0.224940 | 0.224983 | 0.000000 |
| DDF ADDRESS SPACE | 0.005795 | 1.122965 | 0.007127 | 1.135886 | 0.000000 |
| TOTAL | 0.530282 | 1.378663 | 0.269566 | 2.178511 | 0.000000 |

IDUG
Leading the DB2 User
Community since 1988

IDUG DB2 Tech Conference
Prague, Czech Republic | November 2014

#IDUG
33

# DB2 Stats can tell us...

- Worth monitoring LC24 counter as possible indication of zIIP starvation issues from V10 CM onwards
  - LC24 counter covers serialisation latches for EDM thread storage (latch 24) but also Buffer Manager (latch 56)
  - Reported in IFCID 001, QVLSLC24
  - Significant LC24 activity with no corresponding EDM pool activity suggests potential issue with zIIP prefetch



Latch Class 24

IDUG
Leading the DB2 User
Community since 1988

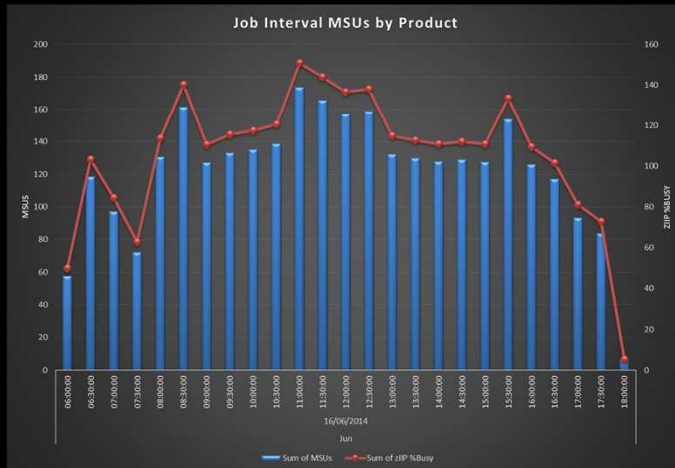IDUG DB2 Tech Conference
Prague, Czech Republic | November 2014

#IDUG
35

## Topics

- Introduction
- What Are zIIPs and How Can They Help?
- How Does DB2 Exploit zIIPs?
- Performance Considerations
- zIIP Usage Measurement and Capacity Planning
- **Case Study**
- Summary

IDUG
Leading the DB2 User
Community since 1988

**IDUG DB2 Tech Conference**
Prague, Czech Republic | November 2014

#IDUG
36

# Case Study – Estimating the Benefit

- Large customer with significant DB2 workload, running on zEC10 hardware
- Single zIIP engine was configured on CEC
- SMF 30 chart shows approx. 120-170 MSU of CP time used by DBM1 during most of the day
  - Contributing to 4HRA
  - zIIP Busy line shows **expected** zIIP usage at V10: 100-160% busy = 2 fully loaded zIIPs!

# Case Study – Quantifying the Benefit

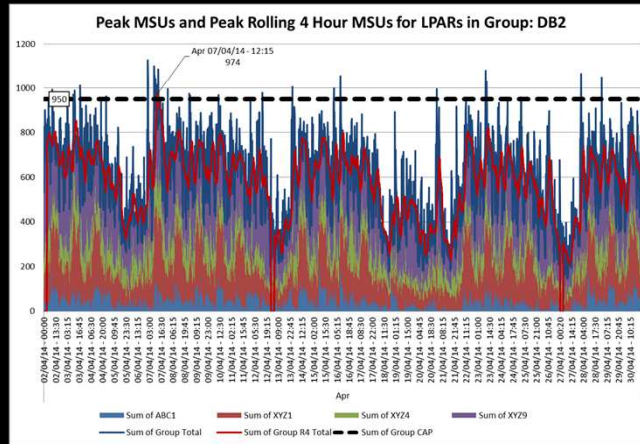- Sub-Capacity Reporting Tool (SCRT) report shows which LPARs contribute to the DB2 MLC

```
==P5=========================================================
PRODUCT MAX CONTRIBUTORS
```

| | | | | LPAR | LPAR | LPAR | LPAR | LPAR | LPAR | LPAR | LPAR |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Product N | Product ID | Highest | Date/Time | ABC1 | ABC2 | ABC3 | XYZ1 | XYZ4 | XYZ8 | XYZ9 | DEV1 |
| z/OS V1 | 5694-A01 | 977 | 07 Apr 2014 - 11:00 U | 100 | 4 | 2 | 344 | 152 | 17 | 353 | 5 |
| DB2 V9 fo | 5635-DB2 | 949 | 07 Apr 2014 - 11:00 | 100 | 0 | 0 | 344 | 152 | 0 | 353 | 0 |
| CICS TS fo | 5655-M15 | 949 | 07 Apr 2014 - 11:00 U | 100 | 0 | 0 | 344 | 152 | 0 | 353 | 0 |
| MQSeries | 5655-L82 | 849 | 07 Apr 2014 - 11:00 U | 0 | 0 | 0 | 344 | 152 | 0 | 353 | 0 |
| Tivoli Net | 5697-ENV | 970 | 07 Apr 2014 - 11:00 U | 100 | 4 | 0 | 344 | 152 | 17 | 353 | 0 |
| IBM Enter | 5655-S71 | 316 | 16 Apr 2014 - 22:00 U | 84 | 0 | 0 | 0 | 232 | 0 | 0 | 0 |

**IDUG**
Leading the DB2 User
Community since 1988

**IDUG DB2 Tech Conference**
Prague, Czech Republic | November 2014

🐦 #IDUG

38

## Case Study – Quantifying the Benefit

- Using SMF70, group the LPARs that contribute to the DB2 MLC to verify that DBM1 peaks correspond to DB2 MLC monthly peaks so we know we'll see savings
- The addition of 2 zIIPs on this machine would allow offload of DBM1 address space and reduce peak 4HRA by 120+ MSUs
- 120 MSUs less for software on those LPARs



Peak MSUs and Peak Rolling 4 Hour MSUs for LPARs in Group: DB2

## Case Study – The Payback

- Other CEC showed very similar workload characteristics, so two additional zIIPs were purchased for each CEC (total of 4 zIIPs) prior to DB2 10 upgrade
- DB2 10 showed the expected zIIP offload following upgrade a few weeks later
- Annual MLC savings estimated at approx. £800k per annum
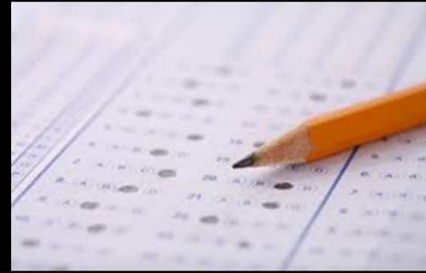- ROI of less than two months

## Topics

- Introduction
- What Are zIIPs and How Can They Help?
- How Does DB2 Exploit zIIPs?
- Performance Considerations
- zIIP Usage Measurement and Capacity Planning
- Case Study
- **Summary**

IDUG

Leading the DB2 User
Community since 1988

IDUG DB2 Tech Conference

Prague, Czech Republic | November 2014

🐦 #IDUG

41

## Summary

- zIIPs are a predominantly a TCO management tool **but**
  - They can improve performance in a CPU-starved environment
  - They can degrade performance if they are over-utilised
- DB2 10 starts directing performance-critical workload to zIIPs for the first time
  - Impact of over-utilised zIIP becomes much bigger
  - Importance of monitoring is correspondingly higher
- You have plenty of instrumentation available to tell you exactly what's happening on your systems
- Work with your z/OS Sysprog / Capacity planner to see if you can save your company big money