

Performance Optimization and Tuning Techniques for IBM Processors, including IBM POWER8

Learn optimization strategies for the new IBM POWER8 processor

Apply strategies to IBM POWER7 and IBM POWER6 processors

Optimize code performance in POWER environments



Brian Hall
Ryan Arnold
Peter Bergner
Wainer dos Santos Moschetta
Robert Enenkel
Pat Haugen
Michael R. Meissner

Alex Mericas
Philipp Oehler
Berni Schiefer
Brian F. Veale
Suresh Warriar
Daniel Zabawa
Adhemerval Zanella

Redbooks



International Technical Support Organization

**Performance Optimization and Tuning Techniques
for IBM Processors, including IBM POWER8**

July 2014

Note: Before using this information and the product it supports, read the information in “Notices” on page ix.

First Edition (July 2014)

This edition pertains to Power Systems servers based on IBM POWER8 processor-based technology. Specific software levels and firmware levels used are noted throughout the text.

© Copyright International Business Machines Corporation 2014. All rights reserved.

Note to U.S. Government Users Restricted Rights -- Use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM Corp.

Contents

Notices	ix
Trademarks	x
Preface	xi
Authors	xi
Now you can become a published author, too!	xv
Comments welcome	xv
Stay connected to IBM Redbooks	xv
Chapter 1. Optimization and tuning on IBM POWER8	1
1.1 Introduction	2
1.2 Outline of this guide	2
1.3 Conventions that are used in this guide	4
1.4 Background	5
1.5 Optimizing performance on POWER8	6
1.5.1 Lightweight tuning and optimization guidelines	7
1.5.2 Deployment guidelines	14
1.5.3 Deep performance optimization guidelines	18
Chapter 2. The POWER8 processor	21
2.1 Introduction to the POWER8 processor	22
2.2 Using POWER8 features	23
2.2.1 Multi-core and multi-thread	23
2.2.2 Multipage size support: Page sizes (4 KB, 64 KB, 16 MB, and 16 GB)	27
2.2.3 Efficient use of cache and memory	28
2.2.4 Transactional memory (TM)	37
2.2.5 Vector Scalar eXtension (VSX)	40
2.2.6 Decimal floating point	42
2.2.7 In-core cryptography and integrity enhancements	42
2.2.8 On-chip accelerators	44
2.2.9 Storage synchronization (sync, lwsync, lwarx, stwcx, and eieio)	44
2.2.10 Fixed-point load and store quadword instructions	46
2.2.11 Instruction fusion	46
2.2.12 Event-based branches (or user-level fast interrupts)	47
2.2.13 Power management and system performance	47
2.3 Related publications	48
Chapter 3. The POWER Hypervisor	51
3.1 Introduction to the POWER8 Hypervisor	52
3.2 POWER8 virtualization	53
3.2.1 Virtual processors	53
3.2.2 Page table sizes for LPARs	57
3.2.3 Placing LPAR resources to attain higher memory affinity	57
3.2.4 Active memory expansion	60
3.2.5 Optimizing Resource Placement: Dynamic Platform Optimizer	61
3.2.6 Partition compatibility mode	61
3.3 Related publications	61
Chapter 4. AIX	63

4.1	Introduction	64
4.2	Using Power features with AIX	64
4.2.1	Multi-core and multi-thread	64
4.2.2	Multipage size support on AIX	74
4.2.3	Efficient use of cache	78
4.2.4	Transactional memory (TM)	81
4.2.5	Vector Scalar eXtension (VSX)	82
4.2.6	Decimal floating point (DFP)	83
4.2.7	On-chip encryption accelerator	85
4.3	AIX operating system-specific optimizations	86
4.3.1	Malloc	86
4.3.2	Pthread tunables	89
4.3.3	pollset	89
4.3.4	File system performance benefits	89
4.3.5	Direct I/O	90
4.3.6	Concurrent I/O (CIO)	90
4.3.7	Asynchronous I/O	90
4.3.8	I/O completion ports	91
4.3.9	shmat versus mmap	91
4.3.10	Large segment tunable aliasing (LSA)	92
4.3.11	64-bit versus 32-bit ABIs	92
4.3.12	Sleep and wake-up primitives (thread_wait and thread_post)	93
4.3.13	Shared versus private loads	94
4.3.14	Workload partitions (WPARs) shared License Program Product (LPP) installs	95
4.4	AIX preferred practices	96
4.4.1	AIX preferred practices that are applicable to all Power Systems generations.	96
4.4.2	AIX preferred practices that are applicable to POWER7 and POWER8	97
4.5	Related publications	98
Chapter 5.	IBM i	101
5.1	Introduction	102
5.2	Using Power features with IBM i	102
5.2.1	Multi-core and multi-thread	102
5.2.2	Multipage size support on IBM i	103
5.2.3	Vector Scalar eXtension (VSX)	103
5.2.4	Decimal floating point	103
5.3	IBM i operating system-specific optimizations	104
5.3.1	IBM i advanced optimization techniques	104
5.3.2	Performance management on IBM i	105
5.4	Related publications	106
Chapter 6.	Linux	107
6.1	Introduction	108
6.2	Using Power features with Linux	108
6.2.1	Multi-core and multi-thread	108
6.2.2	Multipage size support on Linux	113
6.2.3	Efficient use of cache	113
6.2.4	Transactional memory (TM)	113
6.2.5	Vector Scalar eXtension (VSX)	120
6.2.6	Decimal floating point (DFP)	120
6.2.7	Event-based branches	123
6.3	Linux operating system-specific optimizations	123
6.3.1	GCC, toolchain, and IBM Advance Toolchain	124

6.3.2	Tuning and optimizing malloc	127
6.3.3	Large TOC -mcmmodel=medium optimization	131
6.3.4	POWER7 based distro considerations	131
6.3.5	Split-core considerations	132
6.3.6	KVM on Power considerations	132
6.4	Related publications	132
Chapter 7. Compilers and optimization tools for C, C++, and Fortran		135
7.1	Compiler versions and optimization levels	136
7.2	Advanced compiler optimization techniques	137
7.2.1	Common prerequisites	137
7.2.2	XL compiler family	138
7.2.3	GCC compiler family	140
7.3	Capitalizing on POWER8 features with the XL and GCC compilers	142
7.3.1	In-core cryptography	142
7.3.2	Compiler support for VSX	145
7.3.3	Built-in functions for storage synchronization	147
7.3.4	Data Streams Control Register (DSCR) controls	148
7.3.5	Transactional memory (TM)	149
7.4	IBM Feedback Directed Program Restructuring (FDPR)	149
7.4.1	Introduction	149
7.4.2	FDPR supported environments	151
7.4.3	Acceptable input formats	151
7.4.4	General operation	151
7.4.5	Instrumentation and profiling	152
7.4.6	Optimization	154
7.5	Using the Advance Toolchain with IBM XLC and XLF	158
7.6	Related publications	159
Chapter 8. Java		161
8.1	Java levels	162
8.2	32-bit versus 64-bit Java	162
8.3	Memory and page size considerations	163
8.3.1	Medium and large pages for Java heap and code cache	163
8.3.2	Configuring large pages for Java heap and code cache	164
8.3.3	Prefetching	165
8.3.4	Compressed references	165
8.3.5	JIT code cache	166
8.3.6	Shared classes	166
8.3.7	In-core Advanced Encryption Standard (AES) acceleration	167
8.3.8	Transactional memory (TM)	167
8.3.9	Runtime instrumentation	168
8.4	Java garbage collection tuning	168
8.4.1	GC strategy: Optthruput	168
8.4.2	GC strategy: Optavgpause	168
8.4.3	GC strategy: Gencon	169
8.4.4	GC strategy: Balanced	169
8.4.5	Optimal heap size	170
8.5	Application scaling	170
8.5.1	Choosing the correct SMT mode	171
8.5.2	Using resource sets (RSETS)	172
8.5.3	Java lock reservation	174
8.5.4	Java GC threads	174

8.5.5 Java concurrent marking	174
8.6 Related publications	175
Chapter 9. DB2	177
9.1 DB2 and the POWER processor	178
9.2 Taking advantage of the POWER processor	178
9.2.1 Affinitization	178
9.2.2 Page sizes	179
9.2.3 Decimal arithmetics	180
9.2.4 Using SMT priorities for internal lock implementation	180
9.2.5 SIMD	180
9.3 Capitalizing on the compilers and optimization tools for POWER	181
9.3.1 Whole-program analysis and profile-based optimizations	182
9.3.2 Feedback directed program restructuring (FDPR)	182
9.4 Capitalizing on POWER virtualization	182
9.4.1 DB2 virtualization	182
9.4.2 DB2 in an AIX workload partition	183
9.5 Capitalizing on the AIX system libraries	183
9.5.1 Using the thread_post_many API	183
9.5.2 File systems	184
9.6 Capitalizing on performance tooling	185
9.6.1 High-level investigation	185
9.6.2 Low-level investigation	185
9.7 Conclusion	186
9.8 Related publications	186
Chapter 10. WebSphere Application Server	189
10.1 IBM WebSphere on Power Systems	190
10.2 Performance and functional considerations	190
10.2.1 Installation	190
10.2.2 Deployment	190
10.2.3 Performance	191
10.2.4 Performance analysis, problem determination, and diagnostic tests	193
Appendix A. Analyzing malloc usage under AIX	195
Introduction	196
How to collect malloc usage information	197
Appendix B. Performance tooling and empirical performance analysis	199
Introduction	200
Performance advisors	200
Expert system advisors	200
Rational Performance Advisor	205
Power Virtualization Performance (PowerVP)	206
AIX	206
CPU profiling	207
AIX trace-based analysis tools	209
Finding emulation issues	214
hpmstat, hpmcount, and tprof -E	215
Linux	215
Empirical performance analysis using the IBM software development kit (SDK) for PowerLinux	216
Using the IBM SDK for PowerLinux Trace Analyzer	217
High library usage	217

Deeper empirical analysis	218
Java (either AIX or Linux).	222
32-bit or 64-bit JDK	222
Java heap size, and garbage collection (GC) policies and parameters	222
Hot method or routine analysis	223
Locking analysis	229
Thread state analysis	229

Notices

This information was developed for products and services offered in the U.S.A.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not grant you any license to these patents. You can send license inquiries, in writing, to:

IBM Director of Licensing, IBM Corporation, North Castle Drive, Armonk, NY 10504-1785 U.S.A.

The following paragraph does not apply to the United Kingdom or any other country where such provisions are inconsistent with local law: INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some states do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

Any references in this information to non-IBM websites are provided for convenience only and do not in any manner serve as an endorsement of those websites. The materials at those websites are not part of the materials for this IBM product and use of those websites is at your own risk.

IBM may use or distribute any of the information you supply in any way it believes appropriate without incurring any obligation to you.

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 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.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to the names and addresses used by an actual business enterprise is entirely coincidental.


COPYRIGHT LICENSE:

This information contains sample application programs in source language, which illustrate programming techniques on various operating platforms. You may copy, modify, and distribute these sample programs in any form without payment to IBM, for the purposes of developing, using, marketing or distributing application programs conforming to the application programming interface for the operating platform for which the sample programs are written. These examples have not been thoroughly tested under all conditions. IBM, therefore, cannot guarantee or imply reliability, serviceability, or function of these programs.

Trademarks

IBM, the IBM logo, and ibm.com are trademarks or registered trademarks of International Business Machines Corporation in the United States, other countries, or both. These and other IBM trademarked terms are marked on their first occurrence in this information with the appropriate symbol (® or ™), indicating US registered or common law trademarks owned by IBM at the time this information was published. Such trademarks may also be registered or common law trademarks in other countries. A current list of IBM trademarks is available on the Web at <http://www.ibm.com/legal/copytrade.shtml>

The following terms are trademarks of the International Business Machines Corporation in the United States, other countries, or both:

Active Memory™	POWER®	PowerLinux™
AIX®	Power Architecture®	PowerPC®
AIX 5L™	POWER Hypervisor™	PowerVM®
Blue Gene/L®	Power Systems™	PowerVP™
DB2®	Power Systems Software™	PureData™
EnergyScale™	POWER6®	Rational®
FDPR®	POWER6+™	Redbooks®
IBM®	POWER7®	Redbooks (logo)  ®
IBM PureData™	POWER7 Systems™	System z®
Micro-Partitioning®	POWER7+™	Tivoli®
Optim™	POWER8™	WebSphere®

VSR, and the Texas Memory Systems logo are trademarks or registered trademarks of Texas Memory Systems, an IBM Company.

Intel, Intel logo, Intel Inside logo, and Intel Centrino logo are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

Linux is a trademark of Linus Torvalds in the United States, other countries, or both.

LTO, the LTO Logo and the Ultrium logo are trademarks of HP, IBM Corp. and Quantum in the U.S. and other countries.

Windows, and the Windows logo are trademarks of Microsoft Corporation in the United States, other countries, or both.

Java, and all Java-based trademarks and logos are trademarks or registered trademarks of Oracle and/or its affiliates.

UNIX is a registered trademark of The Open Group in the United States and other countries.

Other company, product, or service names may be trademarks or service marks of others.

Preface

This IBM® Redbooks® publication focuses on gathering the correct technical information, and laying out simple guidance for optimizing code performance on IBM POWER8™ systems that run the IBM AIX®, IBM i, or Linux operating systems. There is much straightforward performance optimization that can be performed with a minimum of effort and without extensive previous experience or in-depth knowledge.

The POWER8 processor contains many new and important performance features, such as support for eight hardware threads in each core and support for transactional memory. POWER8 is a strict superset of IBM POWER7+™, and so all of the performance features of POWER7+, such as multiple page sizes, also appear in POWER8. Much of the technical information and guidance for optimizing performance on POWER8 presented in this guide also applies to POWER7+ and earlier processors, except where the guide explicitly indicates that a feature is new in POWER8.

This guide strives to focus on optimizations that tend to be positive across a broad set of IBM POWER® processor chips and systems. Specific guidance is given for the POWER8 processor; however, the general guidance is applicable to the IBM POWER7+, IBM POWER7®, IBM POWER6®, IBM POWER5, and even to earlier processors.

This guide is directed to personnel who are responsible for performing migration and implementation activities on IBM POWER8-based servers. This includes system administrators, system architects, network administrators, information architects, and database administrators (DBAs).

Authors

This book was produced by a team of specialists from around the world, working at the International Technical Support Organization, Poughkeepsie Center.

Brian Hall is the lead analyst for Power performance improvement efforts with the IBM Software Group Hardware Acceleration Laboratory team. He works with many IBM software products to capitalize on the IBM Power Architecture® and develop performance preferred practices for software development and deployment. After joining IBM in 1987, Brian originally worked on the IBM XL C/C++/Fortran compilers and on the just-in-time compiler for IBM Java on Power. He has a Bachelor's degree in Computer Science from Queen's University at Kingston and a Master's degree in Computer Science from the University of Toronto.

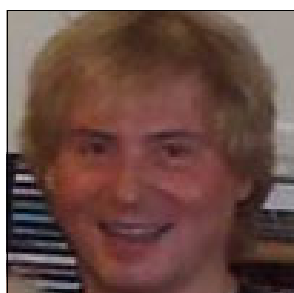
Ryan Arnold is a Senior Software Engineer in the United States. He has 13 years of experience in the computing field. He holds a degree in Computer Science from the University of North Dakota. His areas of expertise include GNU toolchains, system libraries, decimal floating-point math, and ELF ABIs. Ryan has written extensively on ELF ABIs and the GNU C Library.



Peter Bergner is the GCC compiler team lead within the Linux on Power Toolchain department. Since joining IBM in 1996, Peter has worked in a variety of areas, including compiler optimizer development for the IBM i platform, as a core member of the teams that ported Linux and Glibc to 64-bit POWER, and as team lead for the IBM Blue Gene/L® compiler and runtime library development team. He obtained a Ph.D. in Electrical Engineering from the University of Minnesota.



Wainer dos Santos Moschetta is a Staff Software Engineer in the IBM Linux Technology Center, Brazil. He initiated and has led the IBM Software Development Kit (SDK) project for the IBM PowerLinux™ project. He has five years of experience with designing and implementing software development tools for Linux on IBM platforms. Wainer holds a Bachelor's degree in Computer Science from the University of São Paulo. He co-authored Redbooks publication SG24-8075, has published articles and videos for the IBM DeveloperWorks website, and contributes to the IBM PowerLinux technical community blog.



Robert Enenkel has 17 years of experience at the IBM Toronto Lab, both in the optimizing compiler group, and as a research associate at the IBM Centre for Advanced Studies. Robert's interests are in numerical computing as it relates to compilers and operating systems, including floating-point arithmetic, mathematical function libraries, and the performance tuning of algorithms. He obtained his Ph.D. and M.Sc. in numerical analysis from the University of Toronto. (See <https://www.ibm.com/ibm/cas/canada/research/people/robert.jsp>.)

Pat Haugen is an Advisory Software Engineer in Rochester, Minnesota, in the United States. He has 25 years of experience with compilers. He holds a BS degree in Computer Science from the University of North Dakota. Pat's areas of expertise include IBM PowerPC® architecture and code generation, optimization, and tuning.



Michael R. Meissner has worked on compilers since 1979, supporting many different processors and operating systems over the years. In 1988, he started working on the GNU Compiler Collection (GCC). In 2008, Michael joined IBM, to work on GCC, supporting the IBM Cell processor. Later, he switched to work on GCC supporting the PowerPC architecture, including working on adding POWER7 and POWER8 additions to GCC.

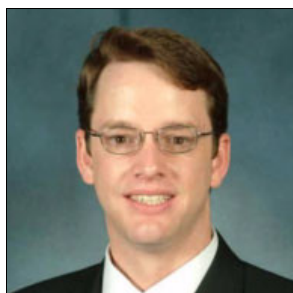
Alex Mericas is a member of the IBM Systems and Technology Group in Austin, Texas. He is a Senior Technical Staff Member and is currently the Performance Architect for POWER8. He designed the performance monitoring unit on POWER4, POWER5, POWER6, POWER7, and PowerPC 970. Alex is an IBM Master Inventor with 47 US patent applications and 22 issued patents covering microprocessor design and hardware performance monitors.



Philipp Oehler is a Microprocessor Development Engineering Professional, working for IBM R&D Germany. He has five years of experience in logic design for cryptographic hardware. He holds a degree in Electrical Engineering (Ph.D.) from the University of Paderborn, Germany, and a degree in Mathematics and Physics (M.Sc.) from the University Innsbruck, Austria. Philipp's areas of expertise include statistical timing analysis and logic synthesis.



Berni Schiefer is a Distinguished Engineer at the IBM Toronto Lab. He is responsible for Information Management performance and benchmarking, specifically for IBM DB2®, IBM PureData™ systems, BigData, MDM, and IBM Optim™ Data Studio performance tools. Berni joined IBM in 1985 and started to work on DB2 in 1991. His current focus is on enhancing the performance and scalability of Information Management solutions. His passion is in introducing advanced technology into products.



Brian F. Veale is a Senior Software Engineer in IBM Power Systems Software™ and an AIX architect. At IBM, he works on the AIX core kernel, specializing in the support and exploitation of new processor features and hardware systems. Brian holds a PhD in Computer Science from the University of Oklahoma, where he was a GAANN Fellow and a lecturer. Prior to obtaining his PhD, Brian worked on training simulators for the U.S. Air Force. He is a Senior member of the Institute of Electrical and Electronics Engineers (IEEE) and a member of the Power Processor Architecture Control Board.



Suresh Warriar is a Senior Technical Staff Member in IBM Power Systems™ Software, specializing in AIX kernel architecture and kernel performance enhancements. Suresh has over 25 years of experience in systems software, including over 15 years with Power Systems, where he leads AIX exploitation of hardware and software technologies. He has a Bachelor's degree in Electrical Engineering from the Indian Institute of Technology, Madras, India, and a Master's Degree in Computer Science from the University of Texas, Austin.

Daniel Zabawa is a compiler developer in Canada with six years of experience in optimization for IBM POWER® Systems. He holds a Masters degree in Computer Science from the University of Toronto. Daniel's areas of expertise include loop optimization and instruction scheduling.



Adhemerval Zanella is a Staff Software Engineer in Brazil. He has nine years of experience with Linux and three years working with toolchain. Adhemerval holds a degree in Computer Science from Universidade de Brasília. His areas of expertise include toolchain, system libraries, and optimization.

Thanks to the following people for their contributions to this project:

- ▶ For peer reviews:
 - Rene Mueller, Research Staff Member, Almaden, California
 - Bruce Semple, Solution Architect, software development, Bethesda, Maryland
 - David Tam, Ph.D., Staff Software Developer, Ontario, Canada
- ▶ Special thanks are in order to a group of subject matter experts (SMEs) who contributed to this book. They provided a perspective that improved the overall value and usability of this book, and they did so with a smile. Many thanks! These persons are:
 - Kent L. Bruinsma, IBM i Performance Tools development (IBM PowerVP™), Rochester, Minnesota
 - Peter J. (Pete) Heyrman, IBM PowerVM® Hypervisor, Rochester, Minnesota
 - Younes Manton, JIT Compiler POWER Optimization, Ontario Canada
 - Bret R. Olszewski, STG Performance, Austin, Texas
 - Robert R. Roediger, Program Optimization, Rochester, Minnesota
 - Todd Rosedahl, Chief Energy Management Engineer on POWER, IBM EnergyScale™ -- Power and Thermal Management, Rochester, Minnesota
- ▶ For technical reviews:
 - Bill Buros, Power Linux Development Architect - LTC, Austin, Texas
 - Chris Conklin, AIX Software Engineer, Austin, Texas
 - Diane Flemming, Power Systems Performance, Austin, Texas
 - Jenifer Hopper, Software Engineer - Linux Performance Analyst, Austin, Texas
 - Steve Munroe, Linux Toolchain Architect and TCEM, Rochester, Minnesota
 - Sergio Reyes, AIX / p-series I/O Subsystem Performance, Austin, Texas
 - Lilian Romero, AIX Commercial Performance, Austin, Texas
 - George Timms, SLIC Supervisor Development, PASE for i, Rochester, Minnesota
 - Scott Vetter, Executive Project Manager, ITSO Redbooks Systems Lead, Austin, Texas
 - Julian Wang, JIT PPC CodeGen, Ontario, Canada
 - Yaakov Yaari, FDP/PR/Linux, Haifa, Israel
- ▶ For overall contributions to this project:
 - International Technical Support Organization, Poughkeepsie Center
 - Ella Buslovich, Graphics Specialist, Poughkeepsie, New York
 - Sertac Cakici, Chef Engineer, Power Systems, Austin, Texas
 - Jessica Erber-Stark, Linux Information Development, Olympia, Washington
 - Karen Lawrence, IBM Redbooks Technical Writer, RTP, North Carolina
 - Yvonne Lyon, IBM Redbooks Editor, ITSO, San Jose, California
 - Rick Peterson, STG Systems Performance, Rochester, Minnesota
 - Tim Pickett, IBM WebSphere® Application Server Performance, Rochester, Minnesota
 - Jessica Rockwood, Senior Manager, DB2 Performance Benchmarks, Ontario, Canada
 - Brian Veale, AIX Kernel Architecture, Austin, Texas

Thanks to the authors of the previous version of this book:

- ▶ Authors of the previous version of this book, *POWER7 and POWER7+ Optimization and Tuning Guide*, SG24-8079, were:
Brian Hall, Mala Anand, Bill Buros, Miso Cilimdžić, Hong Hua, Judy Liu, John MacMillan, Sudhir Maddali, K Madhusudanan, Bruce Mealey, Steve Munroe, Francis P O'Connell, Sergio Reyes, Raul Silvera, Randy Swanberg, Brian Twichell, Brian F Veale, Julian Wang, Yaakov Yaari

Now you can become a published author, too!

Here's an opportunity to spotlight your skills, grow your career, and become a published author—all at the same time! Join an ITSO residency project and help write a book in your area of expertise, while honing your experience using leading-edge technologies. Your efforts will help to increase product acceptance and customer satisfaction, as you expand your network of technical contacts and relationships. Residencies run from two to six weeks in length, and you can participate either in person or as a remote resident working from your home base.

Find out more about the residency program, browse the residency index, and apply online at:

ibm.com/redbooks/residencies.html

Comments welcome

Your comments are important to us!

We want our books to be as helpful as possible. Send us your comments about this book or other IBM Redbooks publications in one of the following ways:

- Use the online **Contact us** review Redbooks form found here:

ibm.com/redbooks

- Send your comments in an email to:

redbooks@us.ibm.com

- Mail your comments to:

IBM Corporation, International Technical Support Organization
Dept. HYTD Mail Station P099
2455 South Road
Poughkeepsie, NY 12601-5400

Stay connected to IBM Redbooks

- Find us on Facebook:

<http://www.facebook.com/IBMRedbooks>

- Follow us on Twitter:

<http://twitter.com/ibmredbooks>

- Look for us on LinkedIn:

<http://www.linkedin.com/groups?home=&gid=2130806>

- Explore new Redbooks publications, residencies, and workshops with the IBM Redbooks weekly newsletter:

<https://www.redbooks.ibm.com/Redbooks.nsf/subscribe?OpenForm>

- Stay current on recent Redbooks publications with RSS Feeds:

<http://www.redbooks.ibm.com/rss.html>



Optimization and tuning on IBM POWER8

This chapter describes the optimization and tuning of IBM POWER8 systems. It covers the following topics:

- ▶ 1.1, “Introduction” on page 2
- ▶ 1.2, “Outline of this guide” on page 2
- ▶ 1.3, “Conventions that are used in this guide” on page 4
- ▶ 1.4, “Background” on page 5
- ▶ 1.5, “Optimizing performance on POWER8” on page 6

1.1 Introduction

The focus of this publication is about gathering the correct technical information, and laying out simple guidance for optimizing code performance on IBM POWER8 systems that run the AIX, IBM i, or Linux operating systems.

This guide strives to focus on optimizations that tend to be positive across a broad set of IBM POWER processor chips and systems. Much of the technical information and guidance for optimizing performance on POWER8 presented in this guide also applies to POWER7+ and earlier processors, except where the guide explicitly indicates that a feature is new in POWER8.

There is much straightforward performance optimization that can be performed with a minimum of effort and without extensive previous experience or in-depth knowledge. This optimization work can provide the following benefits:

- ▶ Substantially improve the performance of the application that is being optimized for POWER8, the focus of this book
- ▶ Typically carry over improvements to systems that are based on related processor chips, such as the IBM POWER7+, IBM POWER7, and IBM POWER6 processor chips
- ▶ Improve performance on other platforms

The POWER8 processor contains many new and important performance features, such as support for eight hardware threads in each core and support for transactional memory. POWER8 is a strict superset of POWER7+, and so all of the performance features of POWER7+, such as multiple page sizes, also appear in POWER8.

This guide is directed to personnel who are responsible for performing migration and implementation activities on IBM POWER8-based servers, which includes systems administrators, system architects, network administrators, information architects, program product developers, software architects, database administrators (DBAs), and compiler writers.

1.2 Outline of this guide

The first chapter of this guide lays out simple strategies for optimizing performance and covers the opportunities that have been found to be the most universally applicable and valuable in past performance efforts (see 1.5, “Optimizing performance on POWER8” on page 6). This chapter is not an exhaustive guide to Power Systems performance, but instead presents a concise overview of typical methodology and areas to focus on in a performance improvement effort. There are references to later chapters in the guide that present a more complete technical discussion of these areas. Later chapters also contain a more complete list of opportunities and techniques to optimize performance that may be valuable in particular cases.

In 1.5.1, “Lightweight tuning and optimization guidelines” on page 7, we describe a set of straightforward steps to set up the environment for performance tuning and optimization, followed by an explanation about how to perform a set of straightforward and easy investigative steps. These steps are the most valuable areas on which to focus in a short performance effort. These steps do not require a deep level of knowledge of the application being optimized, and (with one minor exception) do not involve changing application source code.

In 1.5.2, “Deployment guidelines” on page 14, we describe deployment choices, that is, system setup and configuration choices, so you can tune these designed-for-performance IBM Power Systems for your environment. Together, with section 1.5.1, “Lightweight tuning and optimization guidelines” on page 7, these simple optimization strategies and deployment guidance satisfy the requirements for most environments and can deliver substantial improvements.

Finally, in 1.5.3, “Deep performance optimization guidelines” on page 18, we describe some of the more advanced investigative techniques that can be used to identify performance bottlenecks in an application. It is here that optimization efforts move into the code internals of an application, and improvements are typically made by modifying source code. Of necessity, coverage in this last area is fairly rudimentary, focusing on general areas of investigation and the tooling that you can use.

Most of the remaining material in this guide is technical information that was developed by domain experts at IBM:

- ▶ We provide hardware information about the POWER8 processor (see Chapter 2, “The POWER8 processor” on page 21), highlighting the important features from a performance perspective and laying out the basic information that is drawn upon by the material that follows.
- ▶ Next, we describe the system software stack, examining the IBM POWER Hypervisor™ (see Chapter 3, “The POWER Hypervisor” on page 51), the AIX, IBM i, and Linux operating systems and system libraries (see Chapter 4, “AIX” on page 63, Chapter 5, “IBM i” on page 101, and Chapter 6, “Linux” on page 107), and the compilers (see Chapter 7, “Compilers and optimization tools for C, C++, and Fortran” on page 135). Java (see Chapter 8, “Java” on page 161) also receives extensive coverage.
- ▶ In Chapter 4, “AIX” on page 63, we highlight some of the areas in which AIX exposes some new features of the POWER8 processor. Then, we examine a set of operating system-specific optimization opportunities. The chapter concludes with a short description of AIX preferred practices regarding system setup and maintenance.
- ▶ In Chapter 5, “IBM i” on page 101, we describe IBM i support for a number of features in POWER8 processors (including features available in previous generations of POWER processors). We describe how this operating system can be effective in automatically capitalizing on many new Power architecture features without changes to existing programs. We also provide information about IBM Portable Application Solutions Environment for i (PASE for i), a part of IBM i that allows some AIX application binaries to run on IBM i with little or no changes.
- ▶ In Chapter 6, “Linux” on page 107, we describe the two primary Linux operating systems that are used on POWER8: Red Hat Enterprise Linux (RHEL) for POWER and SUSE Linux Enterprise Server (SLES) for POWER. The minimum supported levels of Linux, including service pack (SP) levels, are RHEL 6.5 (running in POWER7 compatibility mode only) and SLES11 SP3 (running in POWER7 compatibility mode only).

Linux is based on community efforts that are focused not only on the Linux kernel, but also all of the complementary packages, tools, toolchains, and GNU Compiler Collection (GCC) compilers that are needed to effectively use POWER8. IBM provides the expertise for Power Systems by developing, optimizing, and pushing open source changes to the Linux communities.

- ▶ In Chapter 7, “Compilers and optimization tools for C, C++, and Fortran” on page 135, we talk about current compiler versions and optimization levels and how, for projects with increased focus on runtime performance, you can take advantage of the more advanced compiler optimization techniques. We also describe advanced compiler optimization techniques, such as XL compiler programming errors that can provide guidance toward optimization.

- In Chapter 8, “Java” on page 161. we describe the optimization and tuning of Java based applications that are running in a POWER environment.
- Finally, we cover important information about IBM middleware, DB2 (see Chapter 9, “DB2” on page 177) and IBM WebSphere Application Server (see Chapter 10, “WebSphere Application Server” on page 189). Various applications use middleware, and it is critical that the middleware is correctly tuned and performs well. The middleware chapters cover how these products are optimized for POWER8, including select preferred practices for tuning and deploying these products.

The following appendixes are included:

- Appendix A, “Analyzing malloc usage under AIX” on page 195 explains some simple techniques for analyzing how an application is using the system memory allocation routines (*malloc* and related functions in the C library). *malloc* is often a bottleneck for application performance, especially under AIX. AIX has an extensive set of optimized *malloc* implementations, and it is easy to switch between them without rebuilding or changing an application. Knowing how an application uses *malloc* is key to choosing the best memory allocation alternatives that AIX offers. Even Java applications often make extensive use of *malloc*, either in Java Native Interface (JNI) code that is part of the application itself or in the Java class libraries, or in binary code that is part of the software development kit (SDK).
- Appendix B, “Performance tooling and empirical performance analysis” on page 199 describes some of the important performance tools that are available on the IBM Power Architecture under AIX or Linux, and strategies for using them in empirical performance analysis efforts.

These performance tools are most often used as part of the advanced investigative techniques that are described in 1.5, “Optimizing performance on POWER8” on page 6, except for the performance advisors, which are intended as investigative tools that are appropriate for a broader audience of users.

Throughout the book, we include links to related sections among the chapters. For example, Vector Scalar eXtension (VSX) is discussed in the processor chapter, in all of the OS chapters, and in the compiler chapter. Therefore, after the discussion of VSX in the processor chapter, we include links to that same section in the OS chapters and in the compiler chapter.

After you review the advice in this guide, if you would like more information, begin by visiting the IBM Power Systems website:

<http://www.ibm.com/systems/power/index.html>

1.3 Conventions that are used in this guide

In this guide, our convention for indicating sections of code or command examples is shown in Table 1-1.

Table 1-1 Conventions that are used in this guide

Type of example	Format that is used in this guide	Example of our convention
Commands and command options within text	Monofont, bolded	ledit

Type of example	Format that is used in this guide	Example of our convention
Command lines or code examples outside of text	Monofont	<code>ldedit -btextpsize=64k -bdatapsize=64k -bstacksize=64k</code>
Variables in command lines	Monofont, italicized	<code>ldedit -btextpsize=64k -bdatapsize=64k -bstacksize=64k <executable></code>
Variables that are limited to specific choices	Monofont, italicized	<code>-mcmodel={medium large}</code>

1.4 Background

Some continuing trends in processor design are making it more important than ever to invest in analyzing and working to improve application performance. In the past, there were two ways in which newer processor chips delivered higher performance:

- ▶ Increasing the clock rate
- ▶ Making microarchitectural improvements that increase the performance of a single thread

Often, upgrading to a new processor chip gave existing applications a 50% or possibly 100% performance improvement, leaving little incentive to spend much effort to get an uncertain amount of additional performance. However, the approach in the industry has shifted, so that the newer processor chips do not substantially increase clock rates, as compared to the previous generation. In some cases, clock rates declined in newer designs. Recent designs also generally offer more modest improvements in the performance of a single execution thread.

Instead, the focus has shifted to delivering multiple cores per processor chip, and to delivering more hardware threads in each core, known as simultaneous multi-threading (SMT), in IBM Power Architecture terminology. This situation means that some of the best opportunities for improving the performance of an application are in delivering scalable code by having an application make effective use of multiple concurrent threads of execution.

Coupled with the trend toward aggressive multi-core and multi-threaded designs, there are sometimes changes in the amount of cache and memory bandwidth available to each hardware thread. Cache sizes and chip-level bandwidth are, in some cases, increasing at a slower rate than the growth of hardware threads, meaning that the amount of cache per thread is not growing as rapidly. In particular instances, it decreases from one generation to the next. Again, this situation shows where deeper analysis and performance optimization efforts can provide some benefits.

There is also a recent trend toward adding transactional memory support to processors and toward support for special purpose accelerators. Transactional memory is a feature that simplifies multi-threaded programming by providing safe access mechanisms to shared data. Special purpose accelerators may be based on adding new instructions to the core, on chip-level accelerators, or on fast and efficient access mechanisms to new off-chip accelerators, such as graphics processing units (GPUs) or field-programmable gate arrays (FPGAs).

1.5 Optimizing performance on POWER8

This section provides guidance for optimizing performance on POWER8 processor-based systems. We cover the more prominent performance opportunities that have been found in past optimization efforts. The guidance is organized into three broad categories:

1. Lightweight tuning and optimization guidelines:

Lightweight tuning covers simple prescriptive steps for tuning application performance on POWER8. These simple steps can be carried out without detailed knowledge of the internals of the application that is being optimized and usually without modifying the application source code. Simple system utilization and performance tools are used for understanding and improving your application performance. The steps and tools are general guidelines that apply to all types of applications. Although they are simple and straightforward, they often lead to significant performance improvements. It is possible to accomplish these steps in as little as two days or so for a small application, or, at most, two weeks for a large and complex application.

Performance improvement: Consider lightweight tuning to be the starting point for any performance improvement effort.

2. Deployment guidelines:

Deployment guidelines covers tuning considerations that are related to these activities:

- Configuration of a POWER8 system to deliver the best performance
- Associated runtime configuration of the application itself

There are many choices in a deployment, some of which are unrelated to the performance of a particular application, and so, at best, this section can present some guidelines and preferred practices. Understanding logical partitions (LPARs), energy management, I/O configurations, and using multi-threaded cores are examples of typical system considerations that can impact application performance.

Performance improvement: Consider deployment guidelines to be the second required activity for any reasonably extensive performance effort.

3. Deep performance optimization guidelines:

Deep performance analysis covers performance tools and general strategies for identifying and fixing application bottlenecks. This type of analysis requires more familiarity with performance tooling and analysis techniques, sometimes requiring a deeper understanding of the application internals, and often requiring a more dedicated and lengthy effort. Often, a simpler analysis is all that is required to identify serious bottlenecks in an application, but a more detailed investigation is required to perform an exhaustive search for all of the opportunities for increasing performance.

Performance improvement: Consider this the last activity that is undertaken, with simpler analysis steps, for a moderately serious performance effort. The more complex iterative analysis is reserved for only the most performance critical applications.

This chapter provides only minimal background on the guidance provided. Detailed material about these topics is incorporated in the chapters that follow and in the appendixes. The following chapters and appendixes also cover many other performance topics that are not addressed here.

Guidance for POWER8: The guidance that is provided in this book specifically applies to POWER8 processor chips and systems. The guidance that is provided also generally applies to previous generations of POWER processor chips and systems, including POWER7, POWER6, and POWER5. When our guidance is not applicable to all generations of Power Systems, we note that for you.

1.5.1 Lightweight tuning and optimization guidelines

This section covers building and performance testing applications on POWER8, and gives a brief introduction to the most important simple performance tuning opportunities that are identified for POWER8. More details about these and other opportunities are presented in the later chapters of this guide.

Performance test beds and workloads

In performance work, when you are tuning and optimizing an application for a particular processor, you must run and measure performance levels on that processor. Although there are some characteristics that are shared among processor chips in the same family, each generation of processor chip has unique performance features and characteristics. Optimizing code for POWER8 requires that you set up a test bed on a POWER8 system.

You want to see good performance across a range of newer systems, with a special emphasis on optimizing for the latest design. For Power Systems, the previous POWER7 generation is still commonly used, and it may be necessary to support even older POWER6 and earlier systems. For this reason, it is best to have multiple test bed environments, that is, a POWER8 system for most optimization work, and POWER7 and possibly POWER6 systems, for limited testing to ensure that all tuning is beneficial on the previous generations of hardware.

POWER8, POWER7, and POWER6 processors are dissimilar in some respects, and some simple steps can be taken to ensure good performance of a single binary running on any of these systems. In particular, see the information in “C, C++, and Fortran compiler options” on page 9.

Performance test beds must be sized and configured for performance and scalability testing. Choose your scalability goals based on the requirements that are placed on an application, and the test bed must accommodate at least the minimum requirements. For example, when you target a multi-threaded application to scale up to four cores on POWER8, it is important that the test bed be at least a 4-core system and that tests are configured to run in various configurations (1-core, 2-core, and 4-core). You want to be able to measure performance across the different configurations such that the scalability can be computed. Ideally, a 4-core system delivers four times the performance of a 1-core system, but in practice, the scalability is generally less than ideal. Scalability bottlenecks might not be clearly visible if the only testing done for this example were in a 4-core configuration.

With the multi-threaded POWER8 cores (see 2.2, “Using POWER8 features” on page 23), each processor core can be instantiated with one, two, four, or eight logical CPUs within the operating system, so a 4-core server, with SMT8 mode (eight hardware threads per core), means that the operating system is running 32 logical CPUs. Also, larger-core servers are becoming more pervasive, with scaling considerations well beyond 4-core servers.

The performance test bed should be a dedicated LPAR. You must ensure that there is no other activity on the system (including on other LPARs, if any, configured on the system) when performance tests are run. We suggest that initial performance testing be done in a dedicated resource environment to minimize the factors that affect performance. Ensure that the LPAR is running an up-to-date version of the operating system, at the level that is expected for the

typical usage of the application. Keep the test bed in place after any performance effort so that performance can occasionally be monitored, which ensures that later maintenance of an application does not introduce a performance regression.

Choosing the appropriate workloads for performance work is also important. Ideally, a workload has the following characteristics:

- ▶ Be representative of the expected actual usage of the application.
- ▶ Have simple measures of performance that are easily collected and compared, such as run time or transactions/second.
- ▶ Be easy to set up and run in an automated environment, with a fairly short run time for a fast turnaround in performance experiments.
- ▶ Have a low run-to-run variability across duplicated runs, such that extensive tests are not required to obtain a statistically significant measure of performance.
- ▶ Produce a result that is easily tested for correctness.

When an application is being optimized for multiple operating systems, much of the performance work can be undertaken on just one of the operating systems. However, some performance characteristics are operating system-dependent, so some analysis must be performed on each operating system. In particular, perform profiling and lock analysis separately for each operating system to account for differences in system libraries and kernels. Each operating system also has unique scalability considerations. More operating system-specific optimizations are detailed in Chapter 4, “AIX” on page 63, Chapter 6, “Linux” on page 107, and Chapter 5, “IBM i” on page 101.

Build environment and build tools

The build environment, if separate from the performance test bed, must be running an up-to-date operating system. Only recent operating system levels include Application Binary Interface (ABI) extensions to use or control newer hardware features.

Critically, all compilers that are used to build an application need to use up-to-date versions that offer full support for the target processor chip. Older levels of a compiler might tolerate newer processor chips, but they do not capitalize on the unique features of the latest processor chips. For the IBM XL compilers on AIX or Linux, XLC13 and XLF15 are the first compiler versions that have processor-specific tuning for POWER8. For the GCC compiler on Linux, IBM Advance Toolchain version 6.0 contains an updated GCC compiler that is preferred for POWER7, and versions 7.0 and later contain the updates for POWER8. The IBM XL Fortran Compiler is generally preferred over gfortran for the most optimized high floating point performance characteristics.

For the GCC compiler on Linux, the GCC compilers that come with RHEL and SLES recognize and take advantage of the Power Architecture and optimizations. For improved optimizations and newer GCC technology, the IBM Advance Toolchain package provides an updated GCC compiler and optimized toolchain libraries for use with POWER8.

The Advance Toolchain is a key performance technology available for Power Systems running Linux. It includes newer, Power-optimized versions of compilers (GCC, G++, and gfortran), utilities, and libraries, along with various performance tools. The full Advance Toolchain must be installed in the build environment, and the Advance Toolchain runtime package must be installed in the performance test bed. The Toolchain is designed to coexist with the GCC compilers and toolchain that are provided in the standard Linux distributions. More information is available in 6.3.1, “GCC, toolchain, and IBM Advance Toolchain” on page 124.

Along with the compilers for C/C++ and Fortran, there is the separate IBM Feedback Directed Program Restructuring (IBM FDPR®) tool to optimize performance. FDPR takes a post-link executable image (such as one produced by static compilers) and applies additional optimizations. FDPR is another tool that can be considered for optimizing applications that are based on an executable image. More details can be found in 7.4, “IBM Feedback Directed Program Restructuring (FDPR)” on page 149.

Java also contains a dynamic Just-In-Time (JIT) compiler, and only newer versions are tuned for POWER8. However, Java compilations to binary code take place at application execution time, so a newer Java release must be installed on the performance test bed system.

C, C++, and Fortran compiler options

For the static compilers, the important compilation options to consider are as follows:

- ▶ *Basic optimization options:* With the IBM XL compilers, the minimum suggested optimization level to use is **-O**; whereas for GCC it is **-O3**. Higher levels of optimization are better for some types of code, and you might want to experiment with them. The XL compilers also have optimization options, such as **-qhot**, that can be evaluated. More options are detailed in Chapter 7, “Compilers and optimization tools for C, C++, and Fortran” on page 143. The more aggressive optimization options might not work for all programs and might need to be coupled with the strict options described in this list (see “Strict options” on page 9).
- ▶ *Target processor chip options:* It is possible to build a single executable that runs on various POWER processors. However, that executable does not take advantage of some of the features added to later processor chips, such as new instructions. If only a restricted range of newer processor chips must be supported, consider using the compilation options that enable the usage of newer features. With the XL compilers, for example, if the executable must run only on POWER7 or higher processors (including POWER8), the **-qarch=pwr7** option can be specified. The equivalent GCC option is **-mcpu=power7**. Similarly, if the executable must run only on POWER8 processors, the XL **-qarch=pwr8** option can be specified. The equivalent GCC option is **-mcpu=power8**.
- ▶ *Target processor chip tuning options:* The XL compiler **-qtune** option specifies that the code produced must be tuned to run optimally on particular processor chips. The executable that is produced still runs on other processor chips, but might not be tuned for them. The equivalent GCC option is **-mtune**. Here are some possible options to consider:
 - **-qarch=ppc64 -qtune=pwr8** for an executable that is optimized to run on POWER8, but that can run on all 64-bit implementations of the Power Architecture (POWER8, POWER7, POWER6, and so on)
 - **-qarch=pwr7 -qtune=pwr8** for an executable that can run on POWER7 or POWER8 (with access to vector scalar eXtension (VSX) features), but is optimized to run on POWER8
 - **-qarch=pwr5 -qtune=balanced** for an executable that can run on POWER5 and higher chips, and is tuned for good performance for all recent POWER systems (including POWER6, POWER7, POWER8)
 - **-mtune=power7** to tune for POWER7 on GCC and **-mtune=power8** for POWER8 on GCC
- ▶ *Strict options:* Sometimes the compilers can produce faster code by subtly altering the semantics of the original source code. An example of this scenario is expression reorganization. Especially for floating point code, the effect of expression reorganization can produce different results. For some applications, these optimizations must be prevented to achieve valid results. For the XL compilers, certain semantic-altering transformations are allowed by default at higher optimization levels, such as **-O3**, but those transformations can be disabled by using the **-qstrict** option (for example, **-O3**

-qstrict). For GCC, the default is strict mode, but you can use **-ffast-math** to enable optimizations that are not concerned with Not a Number (NaN), signed zeros, infinities, floating point expression reorganization, or setting the `errno` variable. The new **-Ofast** GCC option includes **-O3** and **-ffast-math**, and might include other options in the future.

- ▶ *Source code compatibility options*: The XL compilers assume that the C and C++ source code conforms to language rules for aliasing. On occasion, older source code fails when compiled with optimization, because the code violates the language rules. A workaround for this situation is to use the **-qalias=noansi** option. The GCC workaround is the **-fno-strict-aliasing** option.
- ▶ *Profile Directed Feedback (PDF)*: PDF is an advanced optimization feature of the compilers to consider for performance critical applications.
- ▶ *Interprocedural Analysis (IPA)*: IPA is an advanced optimization feature of the compilers to consider for performance critical applications.

A simple way to experiment with the C, C++, and Fortran compilation options is to repeatedly build an application with different option combinations, and then to run it and measure performance to see the effect. If higher optimization levels produce invalid results, try adding one or both of the **-qstrict** and **-qalias** options with the XL compilers, or **-fno-strict-aliasing** with GCC.

Not all source files must be compiled with the same set of options, but *all* files must be compiled at the minimum optimization level. There are cases where optimization was not used on just one or two important source files, and that caused an application to suffer from substantially reduced performance.

Java options

Many Java applications are performance sensitive to the configuration of the Java heap and garbage collection (GC). Experimentation with different heap sizes and GC policies is an important first optimization step. For generational GC, consider using the options that specify the split between nursery space (also known as the *new* or *young* space) and tenured space (also known as the *old* space). Most Java applications have modest requirements for long-lived objects in the tenured space, but frequently allocate new objects with a short life span in the nursery space.

If 64-bit Java is used, use the **-Xcompressedrefs** option.

By default, newer releases of Java use 64 KB medium pages for the Java heap, which is the equivalent of explicitly specifying the **-X1p64k** option. If older releases of Java are used, we strongly encourage using the **-X1p64k** option. Otherwise, those releases default to using 4 KB pages. Often, there is some additional performance improvement that is seen in using larger 16 MB large pages by using the **-X1p** option. However, using 16 MB pages normally requires explicit configuration by the administrator of the AIX or Linux operating system to reserve a portion of the memory to be used exclusively for large pages. (For more information, see 8.3.2, “Configuring large pages for Java heap and code cache” on page 164.) As such, the medium pages are a better choice for general use, and the large pages can be considered for performance critical applications.

Many Java applications benefit from turning off the default hardware prefetching on POWER8 and POWER7. Some recent Java releases will turn off hardware prefetching by default. For details, see the information in “Tuning to capitalize on hardware performance features” on page 12.

On Power Systems, the **-Xcodecache** option often delivers a small improvement in performance, especially in a large Java application. This option specifies the size of each code cache that is allocated by the JIT compiler for the binary code that is generated for Java

methods. Ideally, all of the compiled Java method binary code fits into a single code cache, eliminating the small penalty that may occur when one Java method calls another method when the binary code for the two methods is in different code caches. To use this option, determine how much code space is being used, and then set the size of the option correctly. The maximum size of each code cache that is allocated is 32 MB, so the largest value that can be used for this option is `-Xcodecache32m`. For more information, see 8.3.5, “JIT code cache” on page 166.

The JIT compiler automatically uses an appropriate optimization level when it compiles Java methods, and recent Java releases automatically fully utilize all of the new features of the target POWER8 processor of the system an application is running on.

For more information about Java performance, see Chapter 8, “Java” on page 161.

Optimized libraries

Optimized libraries are important for application performance. This section covers some considerations that are related to standard libraries for AIX or Linux, libraries for Java, or specialized mathematical subroutine libraries that are available for the Power Architecture.

AIX malloc

The AIX operating system offers various memory allocation packages (the standard `malloc()` and related routines in the C library). The default package offers good space efficiency and performance for single-threaded applications, but it is not a good choice for the scalability of multi-threaded applications. Choosing the correct malloc package on AIX is important for performance. Even Java applications can make extensive use of malloc through JNI code or internally in the Java Runtime Environment (JRE).

Fortunately, AIX offers a number of different memory allocation packages that are appropriate for different scenarios. These different packages are chosen by setting environment variables and do not require any code modification or rebuilding of an application.

Choosing the best malloc package requires some understanding of how an application uses the memory allocation routines. Appendix A, “Analyzing malloc usage under AIX” on page 195 shows how to easily collect the required information. Following the data collection, experiment with various alternatives, alone or in combination. Some alternatives that deliver high performance include these:

- ▶ **Pool malloc:** The pool front end to the malloc subsystem optimizes the allocation of memory blocks of 512 bytes or less. It is common for applications to allocate many small blocks, and pools are particularly space- and time-efficient for that allocation pattern. Thread-specific pools are used for multi-threaded applications. The pool malloc is a good choice for both single-threaded and multi-threaded applications.
- ▶ **Multiheap malloc:** The multiheap malloc package uses up to 32 separate heaps, reducing contention when multiple threads attempt to allocate memory. It is a good choice for multi-threaded applications.

Using the pool front end and multiheap malloc in combination is a good alternative for multi-threaded applications. Small memory block allocations, typically the most common, are handled with high efficiency by the pool front end. Larger allocations are handled with good scalability by the multiheap malloc. A simple example of specifying the pool and multiheap combination is by using the environment variable setting:

```
MALLOCOPTIONS=pool,multiheap
```

For more information malloc alternatives, see Chapter 4, “AIX” on page 63 and “Malloc” on page 86.

Linux Advance Toolchain libraries

The Linux Advance Toolchain contains replacements for various standard system libraries. These replacement libraries are optimized for specific processor chips, including POWER5, POWER6, POWER7, and POWER8. After you install the Linux Advance Toolchain, the dynamic linker automatically has programs to use the library that is optimized for the processor chip type in the system.

The libraries in Linux Advance Toolchain Version 5.0 and later are optimized to use the multi-core facilities in POWER7. Advance Toolchain 7.0 provides optimized libraries for POWER8 (defaults to POWER7).

Mathematical Acceleration Subsystem (MASS) Library and Engineering and Scientific Subroutine Library (ESSL)

The Mathematical Acceleration Subsystem (MASS) libraries contain accelerated scalar, Single Instruction Multiple Data (SIMD), and vector versions of a collection of elementary mathematical functions (such as exp, log, and sin) that run on AIX and Linux. The MASS libraries are included with the XL compilers and are automatically used by the compilers when the `-O3 -qhot=level=1` compilation options are used. The MASS routines can be used automatically with the Advance Toolchain GCC by using the `-mvecLibabi=mass` option, but MASS is not included with GCC and must be separately installed. Explore the use of MASS for applications that use elementary mathematical functions. Substantial performance improvements can occur when you use the vector versions of the functions. The MASS routines do not necessarily provide the same accuracy of results or the same edge-case behavior as standard libraries do.

The Engineering and Scientific Subroutine Library (ESSL) contains an extensive set of advanced mathematical functions and runs on AIX and Linux. Avoid having applications write their own versions of functions, such as the Basic Linear Algebra Subprograms (BLAS). Instead, use the Power optimized versions in ESSL.

java/util/concurrent

For Java, all of the standard class libraries are included with the JRE. One package of interest for scalability optimization is `java/util/concurrent`. Some classes in `java/util/concurrent` are more scalable replacements for older classes, such as `java/util/concurrent/ConcurrentHashMap`, which can be used as a replacement for `java/util/Hashtable`. `ConcurrentHashMap` might be slightly less efficient than `Hashtable` when run in smaller system configurations where scalability is not an issue, so there can be trade-offs. Also, switching packages requires a source code change, albeit a simple one.

Tuning to capitalize on hardware performance features

For almost all applications, using 64 KB pages is beneficial for performance. Newer Linux releases (RHEL5, SLES11, and RHEL6) default to 64 KB pages, and AIX defaults to 4 KB pages. Applications on AIX have 64 KB pages that are enabled through one or a combination of the following methods:

1. Using an environment variable setting:

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

2. Modifying the executable file as follows:

```
ldedit -btextpsize=64k -bdatapsize=64k -bstacksize=64k <executable>
```

3. Using linker options at build time:

```
cc -btextpsize:64k -bdatapsize:64k -bstacksize:64k ...  
ld -btextpsize:64k -bdatapsize:64k -bstacksize:64k ...
```

All of these mechanisms for enabling 64 KB pages can be safely used when the application must run on older hardware or operating system levels that do not support 64 KB pages. When the needed support is not in place, the system simply defaults to using 4 KB pages.

As mentioned in “Java options” on page 10, the newer Java releases default to using 64 KB pages. For Java, it is important that the Java heap space uses 64 KB pages, which are enabled by the **-X1p64k** option in older releases of Java.

Larger 16 MB pages are also supported on the Power Architecture and might provide an additional performance boost when compared to 64 KB pages. However, the usage of 16 MB pages requires explicit configuration by the administrator of the AIX or Linux operating system.

For certain types of non-numerical applications, turning off the default hardware prefetching improves performance. In specific cases, disabling hardware prefetching is beneficial for Java programs, WebSphere Application Server, and DB2. One way to control hardware prefetching is at the partition level, where prefetching is turned off by running the following commands:

- ▶ AIX: **dscrctl -n -s 1**
- ▶ Linux: **ppc64_cpu --dscr=1**

Controlling prefetching in this way might not be appropriate if different applications are running in a partition, because some applications might run best with prefetching enabled. There are also mechanisms to control prefetching at the process level.

Since Java7 SR3 (and Java626 SR4), the JVM defaults to disable hardware prefetch on AIX. A corresponding API on Linux is not supported yet. Option **-XXsetHWPrefetch:os-default** can be used to revert back to the AIX default hardware prefetch setting.

POWER8 allows not only prefetching to be enabled or disabled, but it also allows the fine-tuning of the prefetch engine. Such fine-tuning is especially beneficial for scientific and engineering and memory-intensive applications.¹ Since the effect of hardware prefetching is heavily dependent on the way that an application accesses memory, and also dependent on the cache sizes of a particular Power chip, it is always best to explicitly test the effects of different prefetch settings on each chip the application is expected to run on.

For more information about hardware prefetching and hardware and operating system tuning and usage for optimum performance, see Chapter 2, “The POWER8 processor” on page 21, Chapter 4, “AIX” on page 63, Chapter 6, “Linux” on page 107, and Chapter 5, “IBM i” on page 101.

Scalability considerations

Aside from the scalability considerations already mentioned (such as AIX malloc tuning and java/util/concurrent usage), there is one Linux operating system setting that enhances scalability in some cases: setting `sched_compat_yield` to 1. This task is accomplished by running the following command:

```
sysctl -w kernel.sched_compat_yield=1
```

This setting makes the Completely Fair Scheduler more compatible with earlier versions of Linux. Use this setting for Java environments, such as for WebSphere Application Server. For more information about multiprocessing with the Completely Fair Scheduler, go to this site:

<http://www.ibm.com/developerworks/library/l-completely-fair-scheduler/>

¹ *Making data prefetch smarter: adaptive prefetching on POWER7*, available here:
<http://dl.acm.org/citation.cfm?id=2370837> (available for purchase or with access to the ACM Digital Library)

1.5.2 Deployment guidelines

This section discusses deployment guidelines as they relate to virtualized and non-virtualized environments, and the effect of partition size and affinity on deployments.

Virtualized versus non-virtualized environments

Virtualization is a powerful technique that is applicable to situations where many applications are consolidated onto a single physical server. This consolidation leads to better usage of hardware and simplified system administration. Virtualization is efficient on the Power Architecture, but it does come with some costs. For example, the Virtual I/O Server (VIOS) partition that is allocated for a virtualized deployment consumes a portion of the hardware resources to support the virtualization. For situations where few business-critical applications must be supported on a server, it might be more appropriate to deploy with non-virtualized resources. This situation is particularly true in cases where the applications have considerable network requirements.

Virtualized environments offer many choices for deployment, such as dedicated or non-dedicated processor cores and memory, IBM Micro-Partitioning® that uses fractions of a physical processor core, and memory compression. These alternatives are explored in Chapter 3, “The POWER Hypervisor” on page 51. When you set up a virtualized deployment, it is important that system administrators have a complete understanding of the trade-offs inherent in the different choices and the performance implications of those choices. Some deployment choices, such as enabling memory compression features, can disable other performance features, such as support for 64 KB memory pages.

The POWER8 processor and affinity performance effects

The IBM POWER8 is the latest processor chip in the IBM Power Systems family. The POWER8 processor chip is available in configurations with up to twelve cores per chip, as compared to the POWER7 processor chip, which has up to eight cores per chip. Along with the increased number of cores, the POWER8 processor chip implements SMT8 mode, supporting eight hardware threads per core, as compared to the POWER7, which supported only four hardware threads per core. Each POWER8 processor core supports running in single-threaded mode with one hardware thread, an SMT2 mode with two hardware threads, an SMT4 mode with four hardware threads, or an SMT8 mode with eight hardware threads.

Each SMT hardware thread is represented as a logical processor in AIX, IBM i, or Linux. When the hardware runs in SMT8 mode, the operating system has eight logical processors for each dedicated POWER8 processor core that is assigned to the partition. To gain full benefit from the throughput improvement of SMT, applications must use all of the SMT threads of the processor cores.

Each POWER8 chip has memory controllers that allow direct access to a portion of the memory DIMMs in the system. Any processor core on any chip in the systems can access the memory of the entire system, but it takes longer for an application thread to access the memory that is attached to a remote chip than to access data in the local memory DIMMs.

For more information about the POWER8 hardware, see Chapter 2, “The POWER8 processor” on page 21. The following short description provides some background to help understand two important performance issues that are known as *affinity effects*.

Cache affinity

The hardware threads for each core of a POWER8 processor share a core-specific cache space. For multi-threaded applications where different threads are accessing the same data, it can be advantageous to arrange for those threads to run on the same core. By doing so, the shared data remains resident in the core-specific cache space, as opposed to moving

between different private cache spaces in the system. This enhanced *cache affinity* can provide more efficient utilization of the cache space in the system and reducing the latency of data references.

Similarly, the multiple cores on a POWER8 processor share a chip-specific cache space. Again, arranging the software threads that are sharing the data to run on the same POWER8 processor (when the partition spans multiple chips) often allows more efficient utilization of cache space and reduced data reference latencies.

Memory affinity

By default, the POWER Hypervisor attempts to satisfy the memory requirements of a partition using the local memory DIMMs for the processor cores that are allocated to the partition. For larger partitions, however, the partition might contain a mixture of local and remote memory. For an application that is running on a particular core or chip, the application will run best when using only local memory. This enhanced *memory affinity* reduces the latency of memory accesses.

Partition sizes and affinity

In terms of partition sizes and affinity, this section describes Power dedicated LPARs, shared resource environments, and memory requirements.

Power dedicated LPARs

Dedicated LPAR deployments generally use larger partitions, ranging from just one POWER8 core up to a partition that includes all of the cores and memory in a large symmetric multi-processor (SMP) system. A smaller partition might run a single application, and a larger partition typically runs multiple applications, or multiple instances of a single application. A common example of multiple instances of a single application is in deployments of WebSphere Application Server.

With larger partitions, one of the most important performance considerations is often which cores and memory are allocated to a partition. For partitions of up to the number of cores on the chips used in the system, the POWER Hypervisor attempts to allocate all cores for the partition from a single POWER8 chip and attempts to allocate only memory local to the chip that is used. Those partitions generally and automatically have good cache and memory affinity. However, it might not be possible to obtain resources for each of the LPARs from a single chip.

For example, assume that you have a 32-core system with four chips, each with eight cores. If five partitions are configured, each with six cores, the fifth LPAR will spread across three chips. Start the most important partition first to obtain resources from a single chip. (The order of starting partitions is one consideration in obtaining the best performance for high priority workloads.). This is discussed further in 3.2.3, “Placing LPAR resources to attain higher memory affinity” on page 57.

Another example is when the partition sizes are mixed. Here, starting smaller partitions may consume resources that are spread across many chips, resulting in larger partitions that are spread across multiple chips, which might be contained on a chip if the larger partitions are started first. It is a preferred practice to start higher priority partitions first, so that there is a better opportunity for them to obtain good affinity characteristics in their core and memory allocations. The affinity of the cores and memory that is allocated to a partition can be determined by running the AIX `lsrad -va` command or the Linux `numactl --hardware` command. For more information about partition resource allocation and the `lsrad` command, see Chapter 3, “The POWER Hypervisor” on page 51.

For partitions larger than the number of cores on a chip, the partition always spans more than one chip and has a mixture of local and remote memory. For these larger partitions, it is often useful to manually force good affinity for an application. Manual affinity can be forced by binding applications so that they can run only on particular cores, and by specifying to the operating system that only local memory should be used by the application.

Consider an example where you run four instances of WebSphere Application Server on a partition of 16 cores on a POWER8 system that is running in SMT8 mode. Each instance of WebSphere Application Server is bound to run on four of the cores of the system. Because each of the cores has eight SMT threads, each instance of WebSphere Application Server is bound to 32 logical processors. Note that AIX by default runs a POWER8 system in SMT4 mode and the `smtctl` command can be used to switch the system to SMT8 mode, if needed. Good memory and cache affinity on AIX can therefore be ensured by completing these steps:

1. Set the AIX **MEMORY_AFFINITY** environment variable, typically to the value **MCM**. This setting tells the AIX operating system to use local memory when an application thread requires physical memory to be allocated.
2. Start the four instances of WebSphere Application Server by running the following **execrset** commands, which bind the execution to the specified set of logical processors:
 - **execrset -c 0-31 -m 0 -e** (command to start first WebSphere Application Server instance)
 - **execrset -c 32-63 -m 0 -e** (command to start second WebSphere Application Server instance)
 - **execrset -c 64-95 -m 0 -e** (command to start third WebSphere Application Server instance)
 - **execrset -c 96-127 -m 0 -e** (command to start fourth WebSphere Application Server instance)

Here are some important items to understand in this example:

- ▶ For a particular number of instances and available cores, the most important consideration is that each instance of an application runs only on the cores of one processor chip.
- ▶ Memory and logical processor binding is not done independently because doing so can negatively affect performance.
- ▶ The workload must be evenly distributed over WebSphere Application Server processes for the binding to be effective.
- ▶ There is an assumed mapping of logical processors to cores and chips, which is always established at boot time. This mapping can be altered if the SMT mode of the system is changed by running `smtctl -w now`. It is always best to reboot to change the SMT mode of a partition to ensure that the assumed mapping is in place.

For more information about the **MEMORY_AFFINITY** environment variable, the **execrset** command, and related environment variables and commands, see Chapter 4, “AIX” on page 63.

The same forced affinity can be established on Linux by running **taskset** or **numactl**. Consider these examples:

- ▶ **numactl -C 0-31 -l** <command to start first WebSphere Application Server instance>
- ▶ **numactl -C 32-63 -l** <command to start second WebSphere Application Server instance>
- ▶ **numactl -C 64-95 -l** <command to start third WebSphere Application Server instance>
- ▶ **numactl -C 96-127 -l** <command to start fourth WebSphere Application Server instance>

The **-l** option on these **numactl** commands is the equivalent of the AIX **MEMORY_AFFINITY=MCM** environment variable setting.

Even for partitions contained on a single chip, better cache affinity can be established with multiple application instances by using logical processor binding commands. With partitions contained on a single chip, the performance effects typically range up to about 10% improvement with binding. For partitions that span more than one POWER8 processor chip, using manual affinity results in significantly higher performance. More details about this topic are provided in Chapter 3, “The POWER Hypervisor” on page 51.

Shared resource environments

Virtualized deployments that share cores among a set of partitions also can use logical processor binding to ensure good affinity within the guest operating system. However, the real dispatching of physical cores is handled by the underlying host operating system (POWER Hypervisor).

The POWER Hypervisor uses a three-level affinity mechanism in its scheduler to enforce affinity as much as possible. The reason why absolute affinity is not always possible is that partitions can expand and use unused cycles of other LPARs. This process is done using uncapped mode in Power, where the uncapped cycles might not always have affinity. Therefore, binding logical processors that are seen at the operating system level to physical threads seen at the hypervisor level works only in some cases in shared partitions. Achieving a high level of affinity is difficult when multiple partitions share resources from a single pool, especially at high utilization, and when partitions are expanding to use other partition cycles. Therefore, creating large shared processor core pools that span across chips tends to create remote memory accesses. For this reason, it might be less desirable to use larger partitions and large processor core pools where high-level affinity performance is expected.

Virtualized deployments can use micro-partitioning, where a partition is allocated a fraction of a core. Micro-partitioning allows a core allocation as small as 0.1 cores in older firmware levels, and as small as 0.05 cores in more recent firmware levels, when coupled with supporting operating system levels. This powerful mechanism provides great flexibility in deployments. However, very small core allocations can be more appropriate for situations in which many virtual machines are often idle. Therefore, active 0.05 core LPARs can use those idle cycles.

Also, there is one negative performance effect in deployments with considerably small partitions, in particular, with 0.1 or less cores at high system utilization: Java warm-up times can be greatly increased. In a Java execution, the JIT compiler is producing binary code for Java methods dynamically. Steady-state optimal performance is reached after a portion of the Java methods are compiled to binary code. With considerably small partitions, there might be a long warm-up period before reaching steady-state performance, where a 0.05 core LPAR cannot get additional cycles from other LPARs because the other LPARs are consuming their cycles. Also, if the workload that is running on this small-size LPAR does not need more than 5% of a processor core capacity, then the performance impact is mitigated. More details about this topic are provided in Chapter 3, “The POWER Hypervisor” on page 51.

Memory requirements

For good performance, there needs to be enough physical memory available, so that application data does not need to be frequently paged in and out between memory and disk. The physical memory that is allocated to a partition must be enough to satisfy the requirements of the operating system and the applications that are running on the partition.

Java is sensitive to having enough physical memory available to contain the Java heap because Java applications often have frequent GC cycles where large portions of the Java heap are accessed. If portions of the Java heap are paged out to disk by the operating system because of a lack of physical memory, then GC cycles can cause a large amount of disk activity, which is known as *thrashing*.

1.5.3 Deep performance optimization guidelines

Performance tools for AIX and Linux are described in Appendix B, “Performance tooling and empirical performance analysis” on page 199. A deep performance optimization effort typically uses those tools and follows this general strategy:

1. Gather general information about the execution of an application when it is running on a dedicated POWER8 performance system. These are important statistics to consider:
 - The user and system CPU usage of the application: Ideally, a multi-threaded application generates a high overall CPU usage with most of the CPU time in user code. Too high a system CPU usage is generally a sign of a locking bottleneck in the application. Too low an overall usage usually indicates some type of resource bottleneck, such as network or disk. For low CPU usage, look at the number of runnable threads reported by the operating system, and try to ensure that there are as many runnable threads as there are logical processors in the partition.
 - The network utilization of the application: Networks can be a bottleneck in execution either because of bandwidth or latency issues. Link aggregation techniques are often used to solve networking issues.
 - The disk utilization of the application: High disk I/O issues are increasingly being solved by using solid-state disks (SSDs).

Common operating system tools for gathering this general information include **topas** (AIX), **top** (Linux), **vmstat**, **iostat**, and **netstat**. Detailed CPU usage information is available by running **sar**. This command diagnoses cases where some logical processors are saturated and others are underutilized, an issue that is seen with network interrupt processing on Linux.

2. Collect a time-based profile of the application to see where execution time is concentrated. Here are some possible areas of concern:
 - Particular user routines or Java methods with a high concentration of execution time. This situation is an indication of a poor coding practice or an inefficient algorithm that is being used in the application itself.
 - Particular library routines or Java class library methods with a high concentration of execution time. First, determine whether the hot routine or method is legitimately used to that extent. Look for alternatives or more efficient versions, such as using the optimized libraries in the Linux Advance Toolchain or the vector routines in the MASS library (for more information, see “Mathematical Acceleration Subsystem (MASS) Library and Engineering and Scientific Subroutine Library (ESSL)” on page 12).
 - A concentration of execution time in the pthreads library (see “Java profiling example” on page 224) or in kernel locking (see “Where to use” on page 45) routines. This situation is associated with a locking issue. This locking might ultimately arise at the system level (as seen with malloc locking issues on AIX), or at the application level in Java code (associated with synchronized blocks or methods in Java code). The source of locking issues is not always immediately apparent from a profile. For example, with AIX malloc locking issues, the time that is spent in the malloc and free routines might be quite low, with almost all of the impact appearing in kernel locking routines.

The tools for gathering profiles are **tprof** (AIX) and **OProfile** (Linux) (both tools are described in “Rational Performance Advisor” on page 205). The **curt** tool (see “AIX trace-based analysis tools” on page 209) also provides a breakdown, describing where CPU time is consumed and includes more useful information, such as a system call summary.

3. Where there are indications of a locking issue, collect locking information.

With locking problems, the primary concern is to determine where the locking originates in the application source code. Cases such as AIX malloc locking can be easily solved just by switching to a more scalable memory allocation package through the **MALLOCTYPE** and **MALLOCOPTIONS** environment variables. In this case, examine how malloc is used and consider making changes at the source code level. For example, rather than repeatedly allocating many small blocks of memory by calling malloc for each block, the application can allocate an array of blocks and then internally manage the space.

As mentioned in “java/util/concurrent” on page 12, Java locking issues that are associated with some older classes, such as java/util/Hashtable, can be easily solved by using java/util/concurrent/ConcurrentHashMap.

For Java programs, use *Java Lock Monitor* (see “Java Health Center” on page 223). For non Java programs, use the **splat** tool on AIX (see “AIX trace-based analysis tools” on page 209).

4. For Java, the WAIT tool is a powerful, easy-to-use analysis tool that is based on collecting thread state information.

Using the WAIT tool requires installing and running only a data collection shell. The shell collects various information about the Java program execution, the most important of which is a set of javacore files. The javacore files show the state of all of the threads at the time the file was dumped. The collected data is submitted to an online tool using a web browser, and the tool analyzes the data and displays the results with a GUI. The GUI presents information about thread states and has powerful features to drill down to see call chains.

The WAIT tool results combine many of the features of a time-based profile, a lock monitor, and other tools. For Java programs, the WAIT tool might be one of the first analysis tools to consider because of its versatility and ease of use.

For more information about IBM Whole-system Analysis of Idle Time, which is the browser-based (that is, no-install) WAIT tool, go to this website:

<http://wait.researchlabs.ibm.com>

Guidance about POWER processor chips: The guidance that is provided in this publication generally applies to previous generations of POWER processor chips and systems, including POWER7, POWER6, and POWER5. When our guidance is not applicable to all generations of Power Systems, we note that for you.



The POWER8 processor

This chapter introduces the POWER8 processor and describes some of the technical details and features of this product. It covers the following topics:

- ▶ 2.1, “Introduction to the POWER8 processor” on page 22
- ▶ 2.2, “Using POWER8 features” on page 23
- ▶ 2.3, “Related publications” on page 48

2.1 Introduction to the POWER8 processor

The POWER8 processor is manufactured using the IBM 22 nm Silicon-On-Insulator (SOI) technology. Each chip is 567 mm² and contains 1.2 billion transistors. As shown in Figure 2-1, the chip contains twelve cores, each with its own 512 KB L2 and 8 MB L3 (embedded DRAM) cache, two memory controllers, PCIe Gen3 I/O controllers, and an interconnection system that connects all components within the chip. The interconnect also extends through module and board technology to other POWER8 processors in addition to DDR3 memory and various I/O devices. POWER8 systems utilize memory buffer chips to interface between the POWER8 processor and DDR3 or DDR4 memory. Each buffer chip also includes an L4 Cache to reduce the latency of local memory accesses. The number of memory controllers, memory buffer chips, PCIe lanes, and cores available for use depends upon the particular POWER8 system.

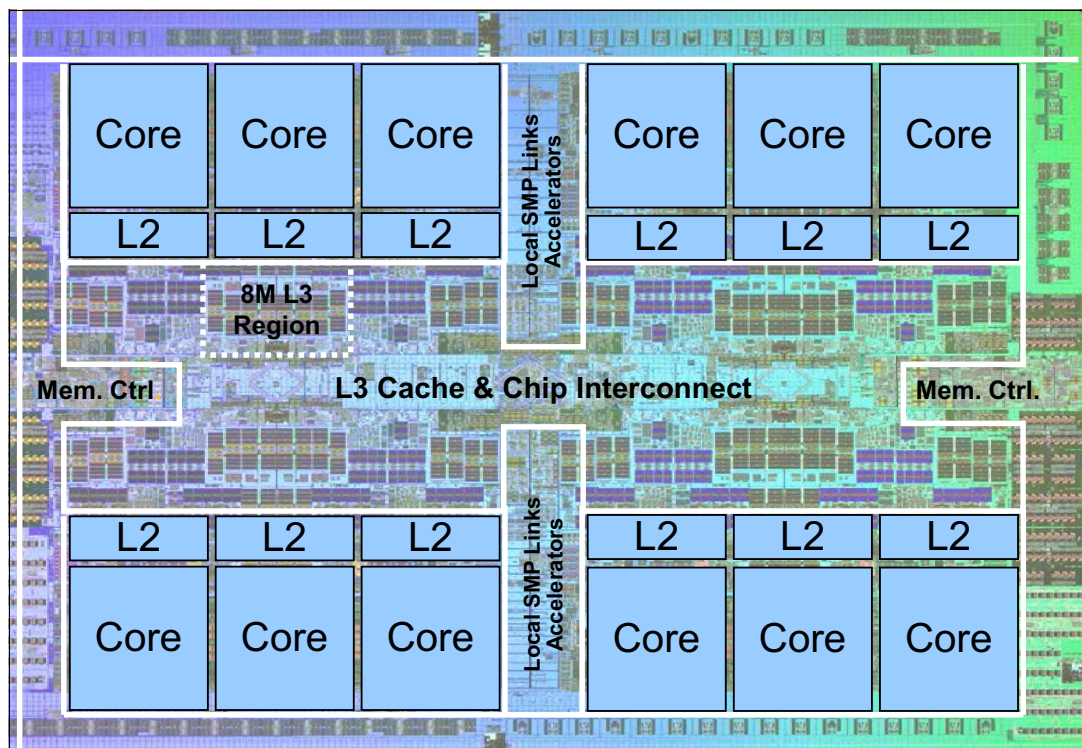


Figure 2-1 The POWER8 processor chip

Each core is a 64-bit implementation of the IBM Power Instruction Set Architecture (ISA) Version 2.07¹ and has the following features:

- ▶ Multi-threaded design, capable of up to eight-way simultaneous multithreading (SMT)
- ▶ 32 KB, eight-way set-associative L1 i-cache
- ▶ 64 KB, eight-way set-associative L1 d-cache
- ▶ 72-entry Effective to Real Address Translation (ERAT) for effective to real address translation for instructions (fully associative)
- ▶ 48-entry primary ERAT (fully associative) and 144-entry secondary ERAT for effective to real address translation for data
- ▶ Aggressive branch prediction, using both local and global prediction tables with a selector table to choose the best predictor

¹ Power ISA Version 2.07, available at <https://www.power.org/documentation/power-isa-version-2-07/>

- ▶ 16-entry link stack
- ▶ 256-entry count cache
- ▶ Aggressive out-of-order execution
- ▶ Two symmetric fixed-point execution units
- ▶ Two symmetric load/store units and two load units, all four of which can also run simple fixed-point instructions
- ▶ An integrated, multi-pipeline vector-scalar floating point unit for running both scalar and SIMD-type instructions, including the Vector Multimedia eXtension (VMX) instruction set and the new Vector Scalar eXtension (VSX) instruction set, and capable of up to eight floating point operations (flops) per cycle (four double precision or eight single precision)
- ▶ In-core Advanced Encryption Standard (AES) encryption capability
- ▶ Hardware data prefetching with 16 independent data streams and software control
- ▶ Hardware decimal floating point (DFP) capability

The POWER8 processor is designed for system offerings from single-socket blades to multi-socket E Class (enterprise) servers. It incorporates a triple-scope broadcast coherence protocol over local and global SMP links to provide superior scaling attributes. Multiple-scope coherence protocols reduce the amount of SMP link bandwidth that is required by attempting operations on a limited scope (single chip or multi-chip group) when possible. If the operation cannot complete coherently, the operation is re-issued using a larger scope to complete the operation.

These are additional features that can augment performance of the POWER8 processor:

- ▶ Adaptive power management
- ▶ Support for DDR3 and DDR4 memory through memory buffer chips that offload the memory support from the POWER8 memory controller
- ▶ 16 MB L4 cache within the memory buffer chip that reduces the memory latency for local access to memory behind the buffer chip; the operation of the L4 cache is transparent to applications running on the POWER8 processor
- ▶ On-chip accelerators, including on-chip encryption, compression, and random number generation accelerators

For more information about this topic, see 2.3, “Related publications” on page 48.

2.2 Using POWER8 features

This section describes several features of POWER8 that can affect performance, including page sizes, cache sharing, SMT priorities, and others.

2.2.1 Multi-core and multi-thread

This section describes the advanced multi-core and multi-thread capabilities of the POWER8 processor. The effective use of the cores and threads is a critically important element of capitalizing on the performance potential of the processor.

Multi-core and multi-thread scalability

POWER8 systems advancements in multi-core and multi-thread scaling are significant. An important POWER8 performance opportunity comes from parallelizing workloads to enable the full potential of the Power platform. Application scaling is influenced by both multi-core and multi-thread technology in POWER8 processors. A single POWER8 chip can contain up to twelve cores. With SMT, each POWER8 core can present eight hardware threads. SMT is the ability of a single physical processor core to simultaneously dispatch instructions from more than one hardware thread context. Because there are multiple hardware threads per physical processor core, additional instructions can run at the same time. SMT is primarily beneficial in commercial environments where the speed of an individual transaction is not as important as the total number of transactions performed. SMT is expected to increase the throughput of workloads with large or frequently changing working sets, such as database servers and web servers.

Additional details about the SMT feature are described in Table 2-1.

Table 2-1 Multi-thread per core features by POWER generation

Technology	Cores/system	Maximum SMT mode	Maximum hardware threads per LPAR
IBM POWER4	32	ST	32
IBM POWER5	64	SMT2	128
IBM POWER6	64	SMT2	128
IBM POWER7	256	SMT4	1024
IBM POWER8	192	SMT8	1536

Information about the multi-thread per core features by single LPAR scaling is available here:

- ▶ Table 4-1 on page 65 (*AIX*)
- ▶ Table 5-1 on page 102 (*IBM i*)
- ▶ Table 6-1 on page 109 (*Linux*)

Operating system enablement usage of multi-core and multi-thread technology varies by operating system and release:

- ▶ Power operating systems present an SMP view of the resources of a partition.
- ▶ Hardware threads are presented as logical CPUs to the application stack.
- ▶ Many applications can use the operating system scheduler to place workloads onto logical processors and maintain the SMP programming model.
- ▶ In some cases, the differentiation between hardware threads per core can be used to improve performance.
- ▶ Placement of a workload on hardware book, drawer and node, socket, core, and thread boundaries can improve application scaling.

Using multi-core and multi-thread features is a challenging prospect.

Further information about this topic, from the OS perspective, is available here:

- ▶ 4.2.1, “Multi-core and multi-thread” on page 64 (*AIX*)
- ▶ 5.2.1, “Multi-core and multi-thread” on page 102 (*IBM i*)
- ▶ 6.2.1, “Multi-core and multi-thread” on page 108 (*Linux*)

For more information about this topic, see 2.3, “Related publications” on page 48.

SMT

The POWER processor architecture uses SMT to provide multiple streams of hardware execution. POWER8 provides eight SMT hardware threads per core and can be configured to run in SMT8, SMT4, SMT2, or single-threaded mode (SMT1 mode or, as referred to in this publication, ST mode). POWER7 and POWER7+ provide four SMT hardware threads per core, and can be configured to run in SMT4, SMT2, or ST mode. POWER6 and POWER5 provide two SMT threads per core, and can be run in SMT2 mode or ST mode.

By using multiple SMT threads, a workload can take advantage of more of the hardware features provided in the Power processor than if a single SMT thread is used per core. By configuring the processor core to run in multi-threaded mode, the operating system can maximize the use of the hardware capabilities that are provided in the system and the overall workload throughput by correctly balancing software threads across all of the cores and SMT hardware threads in the partition.

SMT does include some performance tradeoffs:

- ▶ SMT can provide a significant throughput and capacity improvement on Power processors. When you are in SMT mode, there is a trade-off between overall CPU throughput and the performance of each hardware thread. SMT allows multiple instruction streams to be run simultaneously, but this concurrency can cause some resource conflict between the instruction streams. This conflict can result in a decrease in performance for an individual thread, but an increase in overall throughput.
- ▶ Some workloads do not run well with the SMT feature. This situation is not typical for commercial workloads, but has been observed with scientific (floating point-intensive) workloads.

Information about the topic of SMT, from the OS perspective, is available here:

- ▶ “Simultaneous Multithreading (SMT)” on page 65 (*AIX*)
- ▶ “SMT” on page 102 (*IBM i*)
- ▶ “Simultaneous multithreading (SMT)” on page 109 (*Linux*)

SMT priorities

POWER5 introduced the capability for the SMT thread priority level for each hardware thread to be set, controlling the relative priority of the threads within a single core. This allows each SMT thread to be adjusted so that it can receive more or less favorable performance than the other threads in the same core. The relative difference between the priority of each hardware thread determines the number of decode cycles each thread receives during a period.² This mechanism can be used in various situations, for example, to boost the performance of other threads on the same processor core, while the thread with a lowered priority is waiting on a lock, or waiting on other cooperative threads to reach a synchronization point.

Table 2-2 lists various SMT thread priority levels that are supported in the Power architecture. The level to which the code can set the SMT priority level is controlled by the privilege level that the code is running at (such as problem-state versus supervisor level). For example, code that is running in problem-state cannot set the SMT priority level to High.

Changing the SMT priority level can generally be done in one of the following ways:

- ▶ Executing a Priority Nop, a special form of the `or x,x,x nop`
- ▶ Writing a value to the Program Priority Register (PPR) by executing `mtppr`
- ▶ Through a system call, which can be used by problem-state programs to set priorities in the range permitted for supervisor state

² `thread_set_smt_priority` or `thread_read_smt_priority` System Call, available here:

http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.kerneltechref/doc/ktechrfl/thread_set_smt_priority.htm

On POWER5, POWER6 and POWER7, problem-state programs can only set thread priority values in the range of low (2) to medium (4). POWER7+ additionally introduced allowing a problem-state program to set the thread priority value to very low (1). POWER8 now introduces the ability for a problem-state program to temporarily change the thread priority value to medium-high (5). However, access to medium-high priority is controlled by the operating system through the new Problem State Priority Boost Register introduced in POWER8.³

Table 2-2 SMT thread priority levels for POWER5, POWER6, POWER7, POWER7+, and POWER8^{4, 5}

SMT thread priority level	PPR (11:13)	Priority Nop	Minimum privilege required to set level in POWER5, POWER6, and POWER7	Minimum privilege required to set level in POWER7+	Minimum privilege required to set level in POWER8
Thread shutoff (read only; set by disabling thread)	b'000'		Hypervisor	Hypervisor	Hypervisor
Very low	b'001'	or 31,31,31	Supervisor	Problem-state	Problem-state
Low	b'010'	or 1,1,1	Problem-state	Problem-state	Problem-state
Medium low	b'011'	or 6,6,6	Problem-state	Problem-state	Problem-state
Medium	b'100'	or 2,2,2	Problem-state	Problem-state	Problem-state
Medium high	b'101'	or 5,5,5	Supervisor	Supervisor	Problem-state
High	b'110'	or 3,3,3	Supervisor	Supervisor	Supervisor
Very high	b'111'	or 7,7,7	Hypervisor	Hypervisor	Hypervisor

Information about the topic of SMT priorities, from the OS perspective, is available here:

- “SMT priorities” on page 66 (*ALX*)
- “SMT priorities” on page 110 (*Linux*)

Affinitization and binding to hardware threads

Functionally, it does not matter which core in the partition an application thread is running on, or what physical memory the data it is accessing is on. From a performance standpoint, however, software threads that all access the same data are best placed on the SMT threads of the same core, or on the cores of the same chip. Operating systems may provide facilities to bind applications or specific software threads to run on specific SMT threads or cores.

For information about affinitization and binding, from the OS perspective, see these topics:

- “Affinitization and binding” on page 66 (*ALX*)
- “Affinitization and binding” on page 111 (*Linux*)

³ Power ISA Version 2.07, available at <https://www.power.org/documentation/power-isa-version-2-07/>

⁴ The required privilege to set a particular SMT thread priority level is associated with the physical processor implementation that the LPAR is running on, and not the processor compatible mode. Therefore, setting Very Low SMT priority only requires user level privilege on POWER7+ processors, even when running in IBM POWER6-compatible, IBM POWER6+™-compatible, or POWER7-compatible modes.

⁵ Power ISA Version 2.07, available at <https://www.power.org/documentation/power-isa-version-2-07/>

Hybrid thread and core

The POWER8 processor allows the SMT mode of each core in a partition to be independently controlled by the operating system. Exactly how this facility is presented to the users is dependent on the specific operating system release and version.

These are some of the ways that the operating systems can expose this feature:

- ▶ The ability to set all of the cores in a partition to run in a specific SMT mode, such as to disable SMT and run all of the cores in ST mode.
- ▶ The ability to dynamically alter the SMT mode of specific cores based on load. When only a small number of software threads are ready to run, the operating system can lower the SMT mode of the cores to give each of the software threads the highest possible performance. When a large number of software threads are ready to run, the operating system can use higher SMT modes and maximize the overall throughput of the partition.
- ▶ The ability to specify a fixed asymmetric SMT configuration, where some cores are in high SMT mode and others have SMT mode disabled. This configuration allows critical software threads within a workload to receive an ST performance boost, and allows the remaining threads to benefit from SMT mode. Here are some typical reasons to take advantage of this hybrid mode:
 - For an asymmetric workload, where the performance of one thread serializes an entire workload. For example, one master thread dispatches work to many subordinate threads.
 - For software threads that are critical to a system administrator.

Information about this topic, from the OS perspective, is available here:

- ▶ “Hybrid thread and core” on page 71 (*AIX*)
- ▶ “Hybrid thread and core” on page 112 (*Linux*)

2.2.2 Multipage size support: Page sizes (4 KB, 64 KB, 16 MB, and 16 GB)

The virtual address space of a program is divided into segments. The size of each segment can be either 256 MB or 1 TB on a Power System. The virtual address space can also consist of a mix of these segment sizes. The segments are again divided into units, called *pages*. IBM Power Architecture provides support for multiple virtual memory page sizes, which provides performance benefits to an application because of hardware efficiencies that are associated with larger page sizes.^{6,7}

The POWER5+ and later processor chips support four virtual memory page sizes: 4 KB, 64 KB, 16 MB, and 16 GB. The POWER6 and later processor chip also supports using 64 KB pages inside segments along with a base page size of 4 KB.⁸ The 16 GB pages can be used only within 1 TB segments.

Large pages provide multiple technical advantages:⁹

- ▶ Reduced Page Faults and Translation Lookaside Buffer (TLB) Misses: A single large page that is being constantly referenced remains in memory. This feature eliminates the possibility of several small pages often being swapped out.

⁶ *Power ISA Version 2.07*, available at <https://www.power.org/documentation/power-isa-version-2-07/>

⁷ *What's New in the Server Environment of Power ISA v2.06*, a white paper from Power.org, available here: <https://www.power.org/documentation/whats-new-in-the-server-environment-of-power-isa-v2-06/> (registration required)

⁸ *Multiple page size support*, available here: http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.prftungd/doc/prftungd/multiple_page_size_support.htm

- ▶ **Unhindered Data Prefetching:** A large page enables unhindered data prefetch (which is constrained by page boundaries).
- ▶ **Increased TLB Reach:** This feature saves space in the TLB by holding one translation entry instead of n entries, which increases the amount of memory that can be accessed by an application without incurring hardware translation delays.
- ▶ **Increased ERAT Reach:** The ERAT on Power is a first level and fully associative translation cache that can go directly from effective to real address. Large pages also improve the efficiency and coverage of this translation cache as well.

Large segments (1 TB) also provide reduced Segment Lookaside Buffer (SLB) misses, and increases the reach of the SLB. The SLB is a cache of the most recently used Effective to Virtual Segment translations.

The 16 MB and 16 GB pages are intended only for particularly high performance environments. However, 64 KB pages are considered general purpose, and most workloads benefit from using 64 KB pages rather than 4 KB pages.

Information about this topic, from the OS perspective, is available here:

- ▶ 4.2.2, “Multipage size support on AIX” on page 74
- ▶ 5.2.2, “Multipage size support on IBM i” on page 103
- ▶ 6.2.2, “Multipage size support on Linux” on page 113

2.2.3 Efficient use of cache and memory

Hardware facilities for controlling the efficient use of cache and memory are described in this section.

Cache sharing

Power Systems consist of multiple processor cores and multiple processor chips that share caches and memory in the system. The architecture uses a processor and memory layout that you can use to scale the hardware to many nodes of processor chips and memory. One advantage is that systems can be used for multiple workloads and workloads that are large. However, these characteristics must be carefully weighed in the design, implementation, and evaluation of a workload. Aspects of a program, such as the allocation of data across cores and chips and the layout of data within a data structure, play a key role in maximizing performance, especially when scaling across many processor cores and chips.

Power Systems use a cache-coherent SMP design, in which all of the memory in the system is accessible to all of the processor cores in the system, and all of the cache is coherently maintained:

- ▶ Any processor core on any chip can access the memory of the entire system.
- ▶ Any processor core can access the contents of any core cache, even if it on a different chip.

Processor core access: In both of these cases, the processor core can access only memory or cache that it has authorized access to using normal operating system and Hypervisor memory access permissions and controls.

In POWER8 systems, each chip consists of twelve processor cores, each with on-core L1 instruction and d-caches, an L2 cache, and an L3 cache, as shown in Figure 2-2 on page 30.

⁹ *What's New in the Server Environment of Power ISA v2.06*, a white paper from Power.org, available here: <https://www.power.org/documentation/whats-new-in-the-server-environment-of-power-isa-v2-06/> (registration required)

All of these caches are effectively shared. The L2 cache has a longer access latency than L1, and L3 has a longer access latency than L2. Each chip also has memory controllers, allowing direct access to a portion of the memory DIMMs in the system.¹⁰

Thus, it takes longer for an application thread to access data in cache or memory that is attached to a remote chip than to access data in a local cache or memory. These types of characteristics are often referred to as *affinity performance effects* (see “The POWER8 processor and affinity performance effects” on page 14). In many cases, systems that are built around different processor models have varying characteristics (for example, while L3 is supported, it might not be implemented on some models).

Functionally, it does not matter which core in the system an application thread is running on, or what memory the data it is accessing is on. However, this situation does affect the performance of applications, because accessing a remote memory or cache takes more time than accessing a local memory or cache.¹¹ This situation becomes even more imperative with the capability of modern systems to support massive scaling and the resulting possibility for remote accesses to occur across a large processor interconnection complex.

The effect of these system properties can be observed by application threads, because they often move, sometimes rather frequently, between processor cores. This situation can happen for various reasons, such as a page fault or lock contention that results in the application thread being preempted while it waits for a condition to be satisfied, and then being resumed on a different core. Any application data that is in the cache local to the original core is no longer in the local cache, because the application thread moved and a remote cache access is required.¹² Although modern operating systems, such as AIX, attempt to ensure that cache and memory affinity is retained, this movement does occur, and can result in a loss in performance. For an introduction to the concepts of cache and memory affinity, see “The POWER8 processor and affinity performance effects” on page 14.

The IBM POWER Hypervisor is responsible for these capabilities:

- ▶ Virtualization of processor cores and memory that is presented to the operating system
- ▶ Ensuring that the affinity between the processor cores and memory an LPAR is using is maintained as much as possible

However, it is important for application designers to consider affinity issues in the design of applications, and to carefully assess the impact of application thread and data placement on the cores and the memory that is assigned to the LPAR the application is running in.

Various techniques that are employed at the system level can alleviate the effect of cache sharing. One example is to configure the LPAR so that the amount of memory that is requested for the LPAR is satisfied by the memories that are locally available to processor cores in the system (the memory DIMMs that are attached to the memory controllers for each processor core). Here, it is more likely that the POWER Hypervisor is able to maintain affinity between the processor cores and memory that is assigned to the partition, improving performance¹³.

For more information about LPAR configuration and running the AIX `lssrad-va` command to query the affinity characteristics of a partition, see Chapter 3, “The POWER Hypervisor” on page 51. The equivalent Linux command is `numactl --hardware`.

¹⁰ Of NUMA on POWER7 in IBM i, available here:
http://www.ibm.com/systems/resources/pwrsysperf_P7NUMA.pdf

¹¹ Ibid

¹² Ibid

¹³ Ibid

The rest of this section covers multiple topics that can affect application performance, including the effects of cache geometry, alignment of data, and sensitivity to the scaling of applications to more cores.

Cache geometry

Cache geometry refers to the specific layout of the caches in the system, including their location, interconnection, and sizes. These design details change for every processor chip, even within the Power Architecture. Figure 2-2 shows the layout of a POWER8 chip, including the processor cores, caches, and local memory. Table 2-3 shows the cache sizes and related geometry information for POWER8.

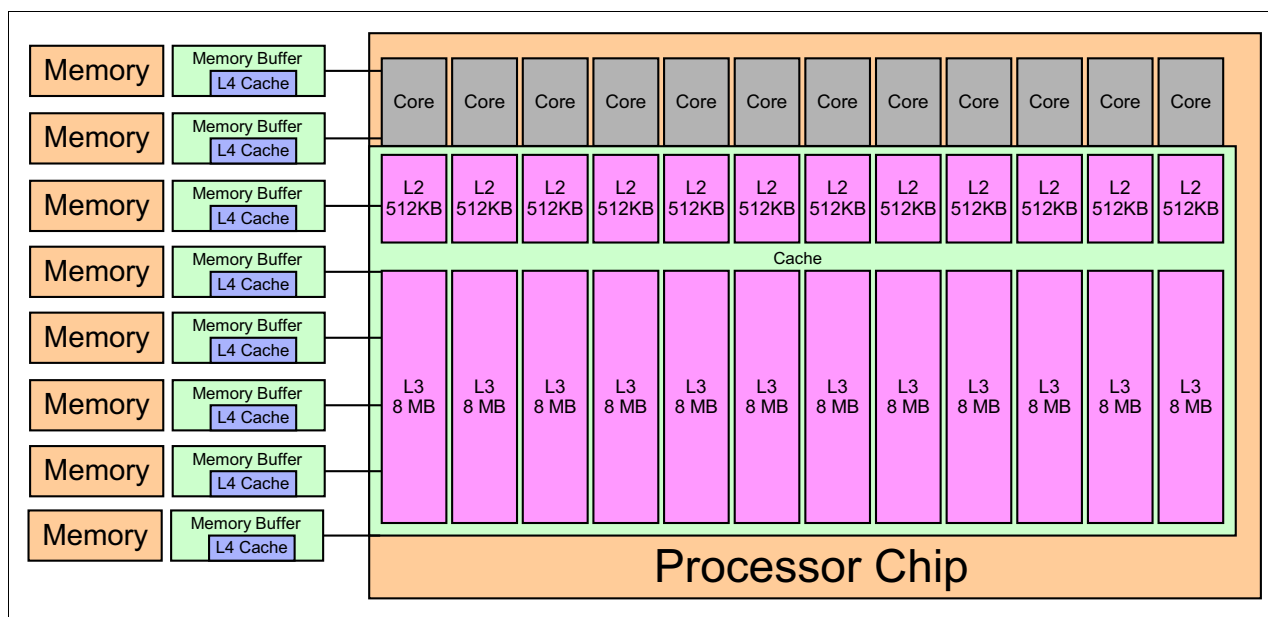


Figure 2-2 POWER8 chip and local memory¹⁴

Table 2-3 POWER8 storage hierarchy¹⁵

Cache	POWER7	POWER7+	POWER8
L1 i-cache: Capacity/associativity	32 KB, 4-way	32 KB, 4-way	32 KB, 8-way
L1 d-cache: Capacity/associativity bandwidth	32 KB, 8-way 2 16 B reads or 1 16 B writes per cycle	32 KB, 8-way 2 16 B reads or 1 16 B writes per cycle	64 KB, 8-way 4 16 B reads or 1 16 B writes per cycle
L2 cache: Capacity/associativity bandwidth	256 KB, 8-way Private 32 B reads and 16 B writes per cycle	256 KB, 8-way Private 32 B reads and 16 B writes per cycle	512 KB, 8-way Private 64 B reads and 16 B writes per cycle
L3 cache: Capacity/associativity bandwidth	On-Chip 4 MB/core, 8-way 16 B reads and 16 B writes per cycle	On-Chip 10 MB/core, 8-way 16 B reads and 16 B writes per cycle	On-Chip 8 MB/core, 8-way 32 B reads and 32 B writes per cycle

¹⁴ Ibid

¹⁵ Ibid

Cache	POWER7	POWER7+	POWER8
L4 cache: Capacity/associativity bandwidth			On-Chip 16MB/buffer chip, 16-way Up to 8 buffer chips per socket

Optimizing for cache geometry

There are several ways to optimize for cache geometry, as described in this section.

Splitting structures into hot and cold elements

A technique for optimizing applications to take advantage of cache is to lay out data structures so that fields that have a high rate of reference (that is, hot) are grouped, and fields that have a relatively low rate of reference (that is, cold) are grouped.¹⁶ The concept is to place the hot elements into the same *byte* region of memory, so that when they are pulled into the cache, they are co-located into the same cache line or lines. Additionally, because hot elements are referenced often, they are likely to stay in the cache. Likewise, the cold elements are in the same area of memory and result in being in the same cache line, so that being written out to main storage and discarded causes less of a performance degradation.

This situation occurs because they have a much lower rate of access. Power Systems use 128-byte length cache lines. Compared to Intel processors (64-byte cache lines), these larger cache lines have the advantage of increasing the reach possible with the same size cache directory, and the efficiency of the cache by covering up to 128-bytes of hot data in a single line. However, it also has the implication of potentially bringing more data into the cache than needed for fine-grained accesses (that is, less than 64 bytes).

As described in *Eliminate False Sharing, Stop your CPU power from invisibly going down the drain*,¹⁷ it is also important to carefully assess the impact of this strategy, especially when applied to systems where there are a high number of CPU cores and a phenomenon referred to as *false sharing* can occur. False sharing occurs when multiple data elements are in the same cache line that can otherwise be accessed independently. For example, if two different hardware threads wanted to update (store) two different words in the same cache line, only one of them at a time can gain exclusive access to the cache line to complete the store. This situation has the following results:

- ▶ Cache line transfers between the processors where those threads are
- ▶ Stalls in other threads that are waiting for the cache line
- ▶ Leaving all but the most recent thread to update the line without a copy in their cache

This effect is compounded as the number of application threads that share the cache line (that is, threads that are using different data in the cache line under contention) is scaled upwards.^{18,17} The discussion about cache sharing¹⁹ in also presents techniques for analyzing false sharing and suggestions for addressing the phenomenon.

Prefetching to avoid cache miss penalties

Prefetching to avoid cache miss penalties is another technique that is used to improve performance of applications. The concept is to prefetch blocks of data to be placed into the cache a number of cycles before the data is needed. This action hides the penalty of waiting for the data to be read from main storage. Prefetching can be speculative when, based on the

¹⁶ *Splitting Data Objects to Increase Cache Utilization (Preliminary Version, 9th October 1998)*. available here: <http://citeseer.uark.edu:8080/citeseerx/viewdoc/summary?doi=10.1.1.84.3359>

¹⁷ *Eliminate False Sharing, Stop your CPU power from invisibly going down the drain*, available here: <http://drdobbs.com/goparallel/article/showArticle.jhtml?articleID=217500206>

¹⁸ Ibid

¹⁹ Ibid

conditional path that is taken through the code, the data might end up not actually being required. The benefit of prefetching depends on how often the prefetched data is used. Although prefetching is not strictly related to cache geometry, it is an important technique.

A caveat to prefetching is that, although it is common for the technique to improve performance for single-thread, single core, and low utilization environments, it actually can decrease performance in high thread-count per-socket and high-utilization environments. Most systems today virtualize processors and the memory that is used by the workload. Because of this situation, the application designer must consider that, although an LPAR might be assigned only a few cores, the overall system likely has a large number of cores. Further, if the LPARs are sharing processor cores, the problem becomes compounded.

The **dcbt** and **dcbtst** instructions are commonly used to prefetch data.^{20,21} *Power Architecture ISA 2.06 Stride N Prefetch Engines to boost Application's performance*²² provides an overview about how these instructions can be used to improve application performance. These instructions can be used directly in hand-tuned assembly language code, or they can be accessed through compiler built-ins or directives.

Prefetching is also automatically done by the POWER8 hardware and is configurable, as described in , “Instruction cache instructions” on page 37.

Alignment of data

Processors are optimized for accessing data elements on their naturally aligned boundaries. Unaligned data accesses might require extra processing time by the processor for individual load or store instructions. They might require a trap and emulation by the host operating system. Ensuring natural data alignment also ensures that individual accesses do not span cache line boundaries.

Similar to the idea of splitting structures into hot and cold elements, the concept of data alignment seeks to optimize cache performance by ensuring that data does not span across multiple cache lines. The cache line size in Power Systems is 128 bytes.

The general technique for alignment is to keep operands (data) on *natural* boundaries, such as a word or doubleword boundary (that is, an int will be aligned on a word boundary in memory). This technique might involve padding and reordering data structures to avoid cases such as the interleaving of chars and doubles: *char; double; char; double*. High-level language compilers are able to ensure optimal data alignment by inserting padding. However, data layout must be carefully analyzed to avoid an undue increase in size by such methods. For example, the previous case of a structure containing *char; double; char; double;* would require 14 bytes of padding. Such an increase in size may result in more cache misses or page misses (especially for rarely referenced groupings of data).

Additionally, to achieve optimal performance, floating point and VMX/VSX have different alignment requirements. For example, the preferred VSX alignment is 16 bytes instead of the element size of the data type being used. This situation means that VSX data that is smaller than 16 bytes in length must be padded out to 16 bytes. The compilers introduce padding as necessary to provide optimal alignment for vector data types.

²⁰ *dcbt (Data Cache Block Touch) instruction*, available here:

http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.aixassem/doc/alangref/idalangref_dcbt_insts.htm

²¹ *dcbtst (Data Cache Block Touch for Store) instruction*, available here:

http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.aixassem/doc/alangref/idalangref_dcbtst_insts.htm

²² *Power Architecture ISA 2.06 Stride N prefetch Engines to boost Application's performance*, available here:

<https://www.power.org/documentation/whitepaper-on-stride-n-prefetch-feature-of-isa-2-06/> (registration required)

Non-vector data intended to be accessed via VSX instructions should be aligned so that VSX loads and stores are performed on addresses aligned to 16-byte boundaries. However, POWER8 improves the handling of misaligned accesses. Most loads that cross cache lines and hit in the d-cache are handled by the hardware with minimal impact on performance.

Byte ordering

The byte ordering (Big-Endian or Little-Endian) is specified by the operating system. In Little Endian mode, byte swapping is performed before data is written to storage and before data is fetched into the execution units. The Load and Store Multiple instructions and the Move Assist instructions are not supported in Little Endian mode. Attempting to execute any of these instructions in Little Endian mode causes the system alignment error handler to be invoked.

POWER8 can operate with the same byte ordering for both instruction and data, or with Split Endian, with instructions and data having different byte ordering.

Sensitivity of scaling to more cores

Different processor chip versions and system models provide less or more scaling of LPARs and workloads to cores. Different processor chips and systems might have different bus widths and latencies. All of these factors result in the sensitivity of the performance of an application/workload to the number of cores it is running on to change based on the processor chip version and system model.

In general terms, an application that tends to not access memory without CPU intervention (that are core-centric) scales perfectly across more cores. Performance loss when scaling across multiple cores tends to come from one or more of the following sources:

- ▶ Increased cache misses (often from invalidations of data by other processor cores, especially for locks)
- ▶ The increased cost of cache misses, which in turn drives overall memory and interconnect fabric traffic into the region of bandwidth limitations (saturating the memory busses and interconnect)
- ▶ The additional cores that are being added to the workload in other nodes, resulting in increased latency in reaching memory and caches in those nodes

Briefly, cache miss requests and returning data can end up being routed through busses that connect multiple chips and memory, which have particular bandwidth and latency characteristics. The goal for scaling across multiple cores, then, is to minimize the change in the potential penalties that are associated with cache misses and data requests as the workload size grows.

It is difficult to assess what strategies are effective for scaling to more cores without considering the complex aspects of a specific application. For example, if all of the cores that the application is running across eventually access all of the data, then it might be wise to interleave data across the processor sockets (which are typically a grouping of processor chips) to optimize them from a memory bus utilization point of view.

However, if the access pattern to data is more localized so that, for most of the data, separate processor cores are accessing it most of the time, the application might obtain better performance if the data is close to the processor core that is accessing that data the most (maintaining affinity between the application thread and the data it is accessing). For the latter case, where the data ought to be close to the processor core that is accessing the data, the AIX `MEMORY_AFFINITY=MCM` environment variable can be set to achieve this behavior. For Linux, the equivalent is the `-1` option on a `numactl` command.

When multiple processor cores are accessing the same data and that data is being held by a lock, resulting in the data line in the cache that is invalidated, programs can suffer. This phenomenon is often referred to as *hot locks*, where a lock is holding data that has a high rate of contention. Hot locks result in cache-to-cache intervention and can easily limit the ability to scale a workload because all updates to the lock are serialized.

Tools such as **splat** (see “AIX trace-based analysis tools” on page 209) can be used to identify hot locks. Additionally, the transactional memory (TM) feature can speed up lock-based programs. Learn more about TM in 2.2.4, “Transactional memory (TM)” on page 37.

Hot locks can be caused by the programmer having lock control access to too large an area of data, which is known as *coarse-grained locking*.²³ In that case, the strategy to effectively deal with a hot lock is to split the lock into a set of fine-grained locks, such that multiple locks, each managing a smaller portion of the data than the original lock, now manage the data for which access is being serialized. Hot locks can also be caused by trying to scale an application to more cores than the original design intended. In that case, using an even finer grain of locking might be possible, or changes can be made in data structures or algorithms, such that lock contention is reduced.

Additionally, the programmer must spend time considering the layout of locks in the cache to ensure that multiple locks, especially hot locks, are not in the same cache line because any updates to the lock itself results in the cache line being invalidated on other processor cores. When possible, pad the locks so that they are in their own distinct cache line.

For more information about this topic, see 2.3, “Related publications” on page 48.

Data prefetching using d-cache instructions and the Data Streams Control Register (DSCR)

The hardware data prefetch mechanism reduces the performance impact that is caused by the latency in retrieving cache lines from higher level caches and from memory. The data prefetch engine of the processor can recognize sequential data access patterns in addition to certain non-sequential (stride-N) patterns and initiate prefetching of d-cache lines from L2 and L3 cache and memory into the L1 d-cache to improve the performance of these storage reference patterns.

The Power ISA architecture also provides cache instructions to supply a hint to prefetch engines for data prefetching to override the automatic stream detection capability of the data prefetcher. Cache instructions, such as **dcbt** and **dcbtst**, allow applications to specify stream direction, prefetch depth, and number of units. These instructions can avoid the starting cost of the automatic stream detection mechanism.

The d-cache instructions **dcbt** (d-cache block touch) and **dcbtst** (d-cache block touch for store) affect the behavior of the prefetched lines. The syntax for the assembly language instructions is as follows:²⁴

```
dcbt RA, RB, TH
dcbtst RA, RB, TH
```

- ▶ *RA* specifies a source general-purpose register for Effective Address (EA) computation.
- ▶ *RB* specifies a source general-purpose register for EA computation.
- ▶ *TH* indicates when a sequence of d-cache blocks might be needed.

²³ *Synchronization & Deadlock Notes*, lecture at Harvard University, available here:

<http://www.read.seas.harvard.edu/~kohler/class/05s-osp/notes/notes8.html>

²⁴ *Power ISA Version 2.07*, available at <https://www.power.org/documentation/power-isa-version-2-07/>

The block that contains the byte addressed by the EA is fetched into the d-cache before the block is needed by the program. The program can later perform loads and stores from the block and might not experience the added delay that is caused by fetching the block into the cache.

The Touch Hint (TH) field is used to provide a hint that the program probably loads or stores to the storage locations specified by the EA and the TH field. The hint is ignored for locations that are caching-inhibited or guarded. The encodings of the TH field depend on the target architecture that is selected with the `-m` flag or the `.machine` assembly language pseudo-op.

The `dcbt` and `dcbtst` instructions provide hints about a sequence of accesses to data elements, or indicate the expected use. Such a sequence is called a *data stream*. The range of values for the TH field describing data streams is 0b01000 - 0b01111. A `dcbt` or `dcbtst` instruction in which TH is set to one of these values is said to be a *data stream variant* of `dcbt` or `dcbtst`.

A data stream to which a program can perform *Load* accesses is said to be a *load data stream*, and is described using the data stream variants of the `dcbt` instruction.

A data stream to which a program can perform *Store* accesses is said to be a *store data stream*, and is described using the data stream variants of the `dcbtst` instruction.

The `dcbt` and `dcbtst` instructions can also be used to provide hints about the transient nature of accesses to data elements. If TH=0b10000, the `dcbt` instruction provides a hint that the program will probably soon load from the block that contains the byte addressed by EA, and that the program's need for the block will be transient (this means the time interval during which the program accesses the block is likely to be short). If TH=0b10001, the `dcbt` instruction provides a hint that the program will probably not access the block that contains the byte addressed by EA for a relatively long period of time.

The contents of the DSCR, a special purpose register, affects how the data prefetcher responds to hardware-detected and software-defined data streams.

The layout of the DSCR register is shown in Table 2-4.

Table 2-4 DSCR register layout (field names are defined following the table)

	SWTE	HWTE	STE	LTE	SWUE	HWUE	UNT CNT	URG	LSD	SNSE	SSE	DPFD
0:38	39	40	41	42	43	44	45:54	55:57	58	59	60	61:63

Here we describe the fields in more detail:

- ▶ 39 Software Transient Enable (SWTE):
New field added in the POWER8 processor. Applies the transient attribute to software-defined streams.
- ▶ 40 Hardware Transient Enable (HWTE):
New field added in the POWER8 processor. Applies the transient attribute to hardware-detected streams.
- ▶ 41 Store Transient Enable (STE):
New field added in the POWER8 processor. Applies the transient attribute to store streams.
- ▶ 42 Load Transient Enable (LTE):
New field added in the POWER8 processor. Applies the transient attribute to load streams.

- ▶ 43 Software Unit count Enable (SWUE):
New field added in the POWER8 processor. Applies the unit count to software-defined streams.
- ▶ 44 Hardware Unit count Enable (HWUE):
New field added in the POWER8 processor. Applies the unit count to hardware-detected streams.
- ▶ 45:54 Unit Count (UNITCNT):
New field added in the POWER8 processor. Number of units in data stream. Streams that exceed this count are terminated.
- ▶ 55:57 Depth Attainment Urgency (URG):
New field added in the POWER7+ processor. This field indicates how quickly the prefetch depth can be reached for hardware-detected streams.
- ▶ Bits 58 Load Stream Disable (LDS):
New field added in the POWER7+ processor Disables hardware detection and initiation of load streams.
- ▶ Bits 59 Stride-N Stream Enable (SNSE):
Enables hardware detection and initiation of load and store streams that have a stride greater than a single cache block. Such load streams are detected when LSD = 0 and such store streams are detected when SSE=1.
- ▶ Bits 60 Store Stream Enable (SSE):
Enables hardware detection and initiation of Store streams.
- ▶ Bits 61:63 Default Prefetch Depth (DPFD):
Supplies a prefetch depth for hardware-detected streams and for software-defined streams for which a depth of zero is specified, or for which **dcbt** or **dcbtst** with TH=1010 is *not* used in their description.
- ▶ Bits 55:57 Depth Attainment Urgency (URG):
This field is a new one added in the POWER7+ processor. This field indicates how quickly the prefetch depth can be reached for hardware-detected streams. Values and their meanings are as follows:
 - 0: Default
 - 1: Not urgent
 - 2: Least urgent
 - 3: Less urgent
 - 4: Medium
 - 5: Urgent
 - 6: More urgent
 - 7: Most urgent

Built-ins for DSCR controls are listed in 7.4.2, “FDPR supported environments” on page 151.

The ability to enable or disable the three types of streams that the hardware can detect (load streams, store streams, or stride-N streams), or to set the default prefetch depth, allows empirical testing of any application. There are no simple rules for determining which settings are optimum overall for a application. The performance of prefetching depends on many different characteristics of the application in addition to the characteristics of the specific system and its configuration. Data prefetches are purely speculative, meaning they can improve performance greatly when the data that is prefetched is, in fact, referenced by the application later, but can also degrade performance by expending bandwidth on cache lines that are not later referenced, or by displacing cache lines that are later referenced by the program.

Similarly, setting DPF to a deeper depth tends to improve performance for data streams that are predominately sourced from memory because the longer the latency to overcome, the deeper the prefetching must be to maximize performance. But deeper prefetching also increases the possibility of stream overshoot, that is, prefetching lines beyond the end of the stream that are not later referenced. Prefetching in multi-core processor implementations has implications for other threads or processes that are sharing cache (in SMT mode) or the same system bandwidth.

For information about modifying the Data Streams Control Register (DSCR) value using the XL compiler family, see 7.4.2, “FDPR supported environments” on page 151.

For information about modifying the DSCR value using the XL compiler family, see 7.3.4, “Data Streams Control Register (DSCR) controls” on page 148.

Instruction cache instructions

The **icbt** instruction provides a hint that the program will probably soon execute code from a storage location and that the cache line containing that code will be loaded into the Level 2 cache. For example, see the following instruction:

```
icbt CT, RA, RB
```

Here we describe it in more detail:

- ▶ *RA* specifies a source general-purpose register for EA computation.
- ▶ *RB* specifies a source general-purpose register for EA computation.
- ▶ *CT* indicates the level of cache the block is to be loaded into. The only supported value for POWER8 is 2.

Information about the efficient use of cache, from the OS perspective, is available here:

- ▶ 4.2.3, “Efficient use of cache” on page 78 (*AIX*)
- ▶ 6.2.3, “Efficient use of cache” on page 113 (*Linux*)
- ▶ 7.3.4, “Data Streams Control Register (DSCR) controls” on page 148 (*compilers*)

2.2.4 Transactional memory (TM)

Transactional memory (TM) is a shared-memory synchronization construct that allows process-threads to perform sequences of storage operations that appear to be atomic to other process-threads and applications. This allows for optimistic execution as a means to take advantage of the inherent parallelism that is found in the latest generation of Power Systems.

One of the main uses of TM is the speed up of lock-based programs by using the speculative execution of lock-based critical sections (CSs), without first acquiring a lock. This allows applications that have not been carefully tuned for performance to take advantage of the benefits of fine-grain locking. The transactional programming model also provides productivity gains when developing lock-based shared memory programs.

Applications can also utilize TM to checkpoint and restore architectural state, independent of the atomic storage access guarantees that are provided by TM.

Using transactional memory

To utilize the TM facility in the most basic form, the process-thread marks the beginning and end of the sequence of storage accesses (namely, the transaction) by using the instructions **tbegin.** and **tend.**, respectively. The **tbegin.** instruction initiates transactional execution, during which the loads and stores appear to occur atomically. The **tend.** instruction terminates transactional execution.

A transaction may either succeed or fail. If a transaction succeeds, it is said to commit, and the transaction appears to have executed as a single atomic unit when viewed by other processors and mechanisms. If a transaction fails, it is as if none of the instructions that were part of the transaction were ever executed. The storage updates that were made since the **tbegin.** instruction was executed are rolled back, and control is transferred to a software failure handler.

It is possible to nest transactions within one another, although the support is using a form of nesting called *flattened nesting*. New transactions that are begun during transactional execution are subsumed by the pre-existing transaction. The effects of a successful nested transaction do not become visible until the outermost (the first transaction that was started in the absence of any previous transactional execution) transaction commits. When a nested transaction fails, the entire set of transactions is rolled back, and control is transferred to the failure handler of the outermost transaction.

A transaction may be put into suspended state by the application, using the **tsuspend.** instruction. This allows a sequence of instructions within the transaction to have the same effect as if the sequence were executed in the absence of a transaction. For example, such instructions are not executed speculatively, and any storage updates will be committed, regardless of transaction success or failure. The **tresume.** instruction is used to resume the transaction and to continue speculative execution of instructions.

Checkpoint state

When a transaction is initiated, and when it is restored following transaction failure, a set of registers is saved or restored, representing the checkpoint state of the processor (for example, the pre-transactional state). The checkpoint state includes all of the problem state, writable registers, with the exception of CR0, FXCC, EBBHR, EBBRR, BESCR, the performance monitor registers, and the TM special purpose registers (SPRs).

Note that the checkpoint state is not directly accessible in either supervisor or problem state. Instead, the checkpoint state is copied into the respective registers when the **treclaim.** instruction is executed. This allows privileged code to save or modify values. The checkpoint state is copied back into the speculative registers (from the respective user-accessible registers) when the new **trechkpt.** instruction is executed.

Transaction failure

A transaction may fail for a variety of reasons, which could be either externally-induced or self-induced. External causes include conflicts with the storage accesses of another process thread (for example, they both access the same storage area and one of the accesses is a store). There are many self-induced causes for a transaction to fail, as in this example:

- ▶ Explicitly aborted using a set of conditional and unconditional abort instructions (for example, various forms of the **tabort.** instruction)
- ▶ Too many nested transactions
- ▶ Too many storage accesses performed in transactional state causing state overflow
- ▶ Execution of certain instructions that are disallowed in transactional state (for example, **slbie**, **dcbi**, etc.)

When a transaction fails, a software failure handler may be invoked. This is accomplished by re-directing control to the instruction following the **tbegin.** of the outermost transaction and setting CR0 to **0b1010**. Therefore, when writing a TM program, the **tbegin.** instruction must always be followed with a conditional branch (for example, **beq**), predicated on bit 2 of CR0. The target of the branch should be the software failure handler that is responsible for handling the transaction failure. For comparison, note that when **tbegin.** is successfully executed at the start of the transaction, CR0 is set to either **0b0000** or **0b0100**.

A transaction failure may be of a transient or a persistent type. Transient failures are typically considered temporary failures, whereas persistent failures indicate that it is unlikely that the transaction will succeed if restarted. The failure handler can retry the transaction or employ a different locking construct or logic path, depending on the nature of the failure. When handling transient type failures, applications may find it useful to keep a count of transient failures and to treat the failure as a persistent type failure on reaching a threshold. If the failure is of persistent type, the expectation is that the applications will fall back to non-transactional logic.

Note that when transaction failure occurs while in a suspended state, failure handling occurs after the transaction is resumed using the **tresume.** instruction.

The software failure handler may identify the cause of the transaction failure by examining bits 0:31 of the Transaction EXception And Summary Register (TEXASR), a special purpose register associated with the TM architecture. In particular, bits 0:6 indicate the failure code, and bit 7 indicates if the failure is persistent and if the transaction will likely fail if attempted again. These bits are copied from the **treclaim.** instruction (privileged code) or the **tabort.** instruction (problem state code) used by software to induce a transaction failure.

The Power Architecture Platform reserves a range of failure codes for use by client operating systems and a separate range for use by a hypervisor, leaving a range of codes free for use by software applications:

- ▶ 0x00 – 0x3F is reserved for use by the OS
- ▶ 0x40 – 0xDF is free for use by problem state (application) code
- ▶ 0xE0 – 0xFF is reserved for use by a hypervisor

Problem state code is limited to using transaction failure codes to the range specified above to provide a failure reason when issuing a **tabort.** instruction.

Sample transaction

Example 2-1 is a sample of assembler code, showing a simple transaction that writes the value in GPR 5 into the address in GPR 4, which is assumed to be shared among multiple threads of execution. If the transaction fails due to a persistent cause, the code falls back to an alternate code path at the label `lock_based_update` (the code for the alternate path is not shown) (based on sample code available from Power.org²⁵).

Example 2-1 A transaction that writes to an address that is shared among multiple execution threads

```
trans_entry:
    tbegin.                # Start transaction
    beq-   failure_hdlr    # Handle transaction failure

# Transaction Body
    stw r5, 0(r4)          # Write to memory pointed to by r4.
    tend.                 # End transaction
    b trans_exit

# Failure Handler
failure_hdlr:
    mfspr r4, TEXASRU      # Read high-order half of TEXASR
    andis. r5, r4, 0x0100  # Is the failure persistent?
    bne lock_based_update  # If persistent, acquire lock and
                           # then perform the write.
```

²⁵ Power ISA Transactional Memory, available here:
<https://www.power.org/documentation/power-isa-transactional-memory/> (registration required).

```

        b trans_entry          # If transient, try again.

# Alternate path for obtaining a lock and performing memory updates
# (non-transactional code path):

lock_based_update:

trans_exit:

```

Information about the topic of transactional memory in AIX and Linux environments is available here:

- ▶ 4.2.4, “Transactional memory (TM)” on page 81 (*AIX*)
- ▶ 6.2.4, “Transactional memory (TM)” on page 113 (*Linux*)
- ▶ 8.3.8, “Transactional memory (TM)” on page 167 (*Java*)

Synchronization mechanisms

In multi-thread programs, synchronization mechanisms are used to guarantee that threads have exclusive access to critical sections. Usually, compare-and-swap (CAS - x86_64) or load-link/store-conditional (LLSC - PowerPC) instructions are used to create locks, a synchronization mechanism. The semantics of locks is this: A running program acquires the lock, executes its CSs in a serialized way (only one thread of execution at a time), and releases the lock.

The serialization of threads due to CSs is a bottleneck to achieve high performance in multi-thread programs. There are some techniques for mitigating or removing such performance issues, for example, non-blocking algorithms, lock-free data, and fine-grained locking.

Lock Elision (LE) is another optimization technique that uses Hardware Transaction Memory (HTM) primitives to avoid lock acquiring. It relies on the behavior of some algorithms that do not have mutually exclusive executions of CS. Some examples might include a hash table insertion where updates can be done in parallel, and locks only needed when the same bucket is accessed at same time.

The LE uses an HTM to first try a transaction on a shared data resource. If it is successful, no locks are required. If the transaction cannot succeed, such as during concurrent access, it falls back to the default locking mechanism.

Information about the topic of transactional memory, from the OS and compiler perspectives, is available here:

- ▶ 4.2.4, “Transactional memory (TM)” on page 81 (*AIX*)
- ▶ 6.2.4, “Transactional memory (TM)” on page 113 (*Linux*)
- ▶ 7.3.5, “Transactional memory (TM)” on page 149 (*XL and GCC compiler families*)
- ▶ 8.3.8, “Transactional memory (TM)” on page 167 (*Java*)

2.2.5 Vector Scalar eXtension (VSX)

Vector Scalar eXtension (VSX) in the Power ISA introduced more support for Vector and Scalar Binary flops conforming to the Institute of Electrical and Electronics Engineers - (IEEE)-754 Standard for Floating Point Arithmetic. The introduction of VSX in to the Power Architecture increases the parallelism by providing SIMD execution functionality for floating point double-precision to improve the performance of the HPC applications.

The following VSX features are provided to increase opportunities for vectorization:

- ▶ A unified register file, a set of Vector-Scalar Registers (VSR™), supporting both scalar and vector operations is provided, eliminating the impact of vector-scalar data transfer through storage.
- ▶ Support for word-aligned storage accesses for both scalar and vector operations is provided.
- ▶ Robust support for IEEE-754 for both vector and scalar flops is provided.
- ▶ Support is provided for symmetric AES instructions that include polynomial multiply to support the Galios Counter Mode (GCM).

A 64-entry Unified Register File is shared across VSX, the Binary floating point unit (BFP), VMX, and the DFP unit. The 32 64-bit Floating Point Registers (FPRs), which are used by the BFP and DFP units, are mapped to registers 0 - 31 of the Vector Scalar Registers. The 32 vector registers (VRs) that are used by the VMX are mapped to registers 32 - 63 of the VSRs,²⁶ as shown in Table 2-5.

Table 2-5 The Unified Register File

FPR0		VSR0
FPR1		VSR1
....		
FPR30		
FPR31		
VR0		
VR1		
..		
..		
VR30		VSR62
VR31		VSR63

VSX supports Double Precision Scalar and Vector Operations and Single Precision Vector Operations. VSX instructions numbering 142 are broadly divided into two categories that can operate on 64 vector scalar registers:^{27, 28, 29, 30, 31}

- ▶ Computational instructions: Addition, subtraction, multiplication, division, extracting the square root, rounding, conversion, comparison, and combinations of these operations
- ▶ Non-computational instructions: Loads/stores, moves, select values, and so on

²⁶ *What's New in the Server Environment of Power ISA v2.06*, a white paper from Power.org, available here: <https://www.power.org/documentation/whats-new-in-the-server-environment-of-power-isa-v2-06/> (registration required)

²⁷ *Support for POWER7 processors*, available here: <http://publib.boulder.ibm.com/infocenter/comphelp/v111v131/index.jsp?topic=/com.ibm.xlc111.aix.doc/getstart/architecture.html>

²⁸ *Vector built-in functions*, available here: http://publib.boulder.ibm.com/infocenter/comphelp/v111v131/index.jsp?topic=/com.ibm.xlc111.aix.doc/compiler_ref/vec_intrin_cpp.html

²⁹ *Initialization of vectors (IBM extension)*, available here: http://publib.boulder.ibm.com/infocenter/comphelp/v111v131/index.jsp?topic=/com.ibm.xlc111.aix.doc/language_ref/vector_init.html

In terms of compiler support for vectors, XLC supports vector processing technologies through language extensions on both AIX and Linux. GCC supports using the VSX engine on Linux. XL and GCC C implement and extend the AltiVec Programming Interface specification.

Information about the topic of VSX, from the OS and compiler perspectives, is available here:

- ▶ 4.2.5, “Vector Scalar eXtension (VSX)” on page 82 (*AIX*)
- ▶ 5.2.3, “Vector Scalar eXtension (VSX)” on page 103 (*IBM i*)
- ▶ 6.2.5, “Vector Scalar eXtension (VSX)” on page 120 (*Linux*)
- ▶ 7.3.2, “Compiler support for VSX” on page 145 (*XL and GCC compiler families*)

2.2.6 Decimal floating point

Decimal (base 10) data is widely used in commercial and financial applications. However, most computer systems have only binary (base two) arithmetic. There are two binary number systems in computers, integer (fixed-point) and floating point. Unfortunately, decimal calculations cannot be directly implemented with binary floating point. For example, the value 0.1 needs an infinitely recurring binary fraction, whereas a decimal number system can represent it exactly, as one tenth. So, using binary floating point cannot ensure that results are the same as those results using decimal arithmetic.

In general, DFP operations are emulated with binary fixed-point integers. Decimal numbers are traditionally held in a binary-coded decimal (BCD) format. Although BCD provides sufficient accuracy for decimal calculation, it imposes a heavy cost in performance, because it is usually implemented in software.

IBM POWER6, POWER7, and POWER8 processor-based systems provide hardware support for DFP arithmetic. The POWER6, POWER7, and POWER8 microprocessor cores include a DFP unit that provides acceleration for the DFP arithmetic. The IBM Power instruction set is expanded. 54 new instructions were added to support the DFP unit architecture. DFP can provide a performance boost for applications that are using BCD calculations.³²

Information about this topic, from the OS perspective, is available here:

- ▶ 4.2.6, “Decimal floating point (DFP)” on page 83 (*AIX*)
- ▶ 5.2.4, “Decimal floating point” on page 103 (*IBM i*)
- ▶ 6.2.6, “Decimal floating point (DFP)” on page 120 (*Linux*)

2.2.7 In-core cryptography and integrity enhancements

POWER8 in-core enhancements are targeting applications by the use of symmetric cryptography, Advanced Encryption Standard (AES); security, Secure Hash Algorithms (SHA-2); and cyclic redundancy check (CRC) algorithms. In cryptography, the information is scrambled so that only an authorized receiver can read the message. Asymmetric-key algorithms require two separate keys, one private and one public, whereas symmetric-key algorithms use the same key for encryption and decryption (for example, AES).

³⁰ *Engineering and Scientific Subroutine Library (ESSL)*, available here:

<http://publib.boulder.ibm.com/infocenter/clresctr/vxrx/index.jsp?topic=/com.ibm.cluster.essl.doc/esslbooks.html>

³¹ *What's New in the Server Environment of Power ISA v2.06*, a white paper from Power.org, available here:

<https://www.power.org/documentation/whats-new-in-the-server-environment-of-power-isa-v2-06/> (registration required)

³² *How to Leverage Decimal Floating-Point unit on POWER6 for Linux*, available here:

<https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/Welc%20to%20High%20Performance%20Computing%20%28HPC%29%20Central/page/How%20to%20Leverage%20Decimal%20Floating-Point%20unit%20on%20POWER6%20for%20Linux>

Many applications not only require information protection (for confidentiality) but they also need to ensure that data is not changed when sent to the receiver (for integrity). This is realized by cryptographic hash functions, which take an arbitrary block of data (often called a *message*) and return a fixed-size bit string (called a *message digest* or *digest*). Well-established algorithms are SHA and CRC.

AES

AES was established for the encryption of electronic data by the U.S. National Institute of Standards and Technology (NIST) in 2001 (FIPS PUB 197). AES is a symmetric-key algorithm which processes data blocks of 128 bits (a block cipher algorithm), and, therefore, naturally fits into the 128-bit VSX data flow. The AES algorithm is completely covered in five new instructions, available in *Power ISA Version 2.07*.³³

AES special mode of operation: Galois Counter Mode (GCM)

The AES GCM mode of operation is designed to provide both confidentiality and integrity (for authentication). GCM is defined for block ciphers (block sizes of 128, 192, and 256 bits). The key feature is that Galois Field multiplication (used for authentication) can be computed in parallel, resulting in higher throughput than the authentication algorithms that use chaining modes.

SHA-2

SHA-2 was designed by the U.S. National Security Agency (NSA) and published in 2001 by the NIST (FIPS PUB 180-2). It is a set of four hash functions (SHA-224, SHA-256, SHA-384, and SHA-512) with message digests that are 224, 256, 384, and 512 bits. The SHA-2 functions compute the digest based on 32-bit words (SHA-224 and SHA-256) or 64-bit words (SHA-384 and SHA-512). Different combinations of *rotate* and *xor* vector instructions have been identified to be merged into a new instruction to accelerate the SHA-2 family. The new instruction comes in two varieties:

- ▶ In word (32-bit), targeting SHA-224 and SHA-256
- ▶ In double-word (64 bit), accelerating SHA-384 and SHA-512 (Power ISA v2.07)

CRC

CRC can be seen as an error-detecting code. It is used in storage devices and digital networks to protect data from accidental (or hacker-intended) changes to raw data. Data to be stored or information sent over the network (in a stream) gets a short **checksum** attached (based on the remainder of the polynomial division and modulo operations). CRC is a reversible function, which makes it unsuitable for use in digital signatures, but it is in use for error detection when data is transferred, for example, in an Ethernet network protocol.

CRC algorithms are defined by the different generator polynomial used. For example, an n-bit CRC is defined by an n-bit polynomial. Examples for applications using CRC-32 are Ethernet (Open Systems Interconnection (OSI) physical layer), Serial Advance Technology Attachment (Serial ATA), Moving Picture Experts Group (MPEG-2), GNU Project file compression software (Gzip), and Portable Network Graphics (PNG, fixed 32-bit polynomial). In contrast, Internet Small Computer System Interface (iSCSI) and the Stream Control Transmission Protocol (SCTP transport layer protocol) are based on a different, 32-bit polynomial.³⁴ The enhancements on POWER8 not only focus on a specific application that supports only one single generator polynomial, but they help to accelerate any kind of CRC size, ranging from 8-bit CRC, 16-bit CRC, and 32-bit CRC, to 64-bit CRC.

³³ *Power ISA Version 2.07*, available at <https://www.power.org/documentation/power-isa-version-2-07/>

³⁴ *Optimization of cyclic redundancy-check codes with 24 and 32 parity bits*, Castagnoli, G.; Bräuer, S.; Herrmann, M. (June 1993). IEEE Transactions on Communications 41 (6), available at <http://ieeexplore.ieee.org/xpl/login.jsp?tp=&arnumber=231911&url=http%3A%2F%2Fieeexplore.ieee.org%2Fiel1%2F26%2F5993%2F00231911>

Information about the topic of in-core cryptography, from the OS and compiler perspectives, is available here:

- ▶ 4.2.7, “On-chip encryption accelerator” on page 85 (*AIX*)
- ▶ 7.3.1, “In-core cryptography” on page 142 (*XL and GCC compiler families*)

2.2.8 On-chip accelerators

On-chip accelerators, initially available in POWER7, provide the following benefits:

- ▶ *On-chip encryption*: AIX transparently uses on-chip encryption accelerators. There are no application visible changes or awareness required.
- ▶ *On-chip compression*:
- ▶ *On-chip random number generator*: AIX capitalizes on the on-chip random number generator, providing the advantages of stronger hardware-based random numbers. In some instances, there can also be a performance advantage.

More details about this topic, from the AIX perspective, are available here:

- ▶ 4.2.7, “On-chip encryption accelerator” on page 85 (*AIX*)
- ▶ “AIX /dev/random (random number generation)” on page 86 (*AIX*)

2.2.9 Storage synchronization (**sync**, **lwsync**, **lwarx**, **stwcx**, and **eieio**)

The Power Architecture storage model provides for out-of-order storage accesses, providing opportunities for performance enhancement when accesses do not need to be in order. However, when accessing storage shared by multiple processor cores or shared with I/O devices, it is important that accesses occur in the correct order that is required for the sharing mechanisms that is used.

The architecture provides mechanisms for synchronization of such storage accesses and defines an architectural model that ought to be adhered to by software. Several synchronization instructions are provided by the architecture, such as **sync**, **lwsync**, **lwarx**, **stwcx**, and **eieio**. There are also operating system-specific locking services provided that enforce such synchronization. Software must be carefully designed when you use these mechanisms to ensure optimal performance while providing appropriate data consistency because of their inherent heavyweight nature.

Concepts and benefits

The Power Architecture defines a storage model that provides weak ordering of storage accesses. The order in which memory accesses are performed might differ from the program order and the order in which the instructions that cause the accesses are run.³⁵

The Power Architecture provides a set of instructions that enforce storage access synchronization, and the AIX kernel provides a set of kernel services that provide locking mechanisms and associated synchronization support.³⁶ However, such mechanisms come with an inherent cost because of the nature of synchronization. Thus, it is important to intelligently use the correct storage mechanisms for the various types of storage access scenarios to ensure that accesses are performed in program order while minimizing their impact.

³⁵ *PowerPC storage model and AIX programming: What AIX programmers need to know about how their software accesses shared storage*, by Lyons, *et al*, available here:

<http://www.ibm.com/developerworks/systems/articles/powerpc.html>

³⁶ Ibid

Associated instructions

The following instructions provide various storage synchronization mechanisms:

sync	This instruction provides an ordering function, so that all instructions issued before the sync complete and no subsequent instructions are issued until after the sync completes. ³⁷
lwsync	This instruction provides an ordering function similar to sync , but it is only applicable to load , store , and dcbz instructions that are run by the processor (hardware thread) running the lwsync instruction, and only for specific combinations of storage control attributes. ³⁸
lwarx	This instruction reserves a storage location for subsequent store using a stcxw instruction and notifies the memory coherence mechanism of the reservation. ³⁹
stcxw	This instruction performs a store to the target location only if the location specified by a previous lwarx instruction is not used for storage by another processor (hardware thread) or mechanism, which invalidates the reservation. ⁴⁰
eieio	This instruction creates a memory barrier that provides an order for storage accesses caused by load , store , dcbz , eciwx , and ecowx instructions. ⁴¹
makeitso	New in POWER8, this instruction allows data to push out to the coherence point as quickly as possible. An attempt to execute the makeitso instruction will provide a hint that preceding stores will be made visible with higher priority.
lbarx/stcbx.	These instructions were added in POWER8 and are similar to lwarx/stcxw. , except that they load and store a byte.
lharx/stchx.	These instructions were added in POWER8 and are similar to lwarx/stcxw. , except that they load and store a 16-bit half word.
ldarx/stcdx.	These instructions are similar to lwarx/stcxw. , except that they load and store a 64-bit double word (requires 64-bit mode).
lqarx/stcq.	These instructions were added in POWER8 and are similar to lwarx/stcxw. , except that they load and store a 128-bit quad word (requires 64-bit mode).

Where to use

Care must be taken when you use synchronization mechanisms in any processor architecture because the associated load and store instructions have a heavier weight than normal loads and stores, and the barrier operations have a cost that is associated with them. Thus, it is imperative that the programmer carefully consider when and where to use such operations, so that data consistency is ensured without adversely affecting the performance of the software and the overall system.

³⁷ *sync* (Synchronize) or *dcs* (Data Cache Synchronize) instruction, available here:

http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=%2Fcom.ibm.aix.aixassem%2Fdoc%2Falangref%2Fidalangref_sync_dcs_instrs.htm

³⁸ *PowerPC storage model and AIX programming: What AIX programmers need to know about how their software accesses shared storage*, Michael Lyons, et al, available here:

<http://www.ibm.com/developerworks/systems/articles/powerpc.html>

³⁹ *lwarx* (Load Word and Reserve Indexed) instruction, available here:

http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=%2Fcom.ibm.aix.aixassem%2Fdoc%2Falangref%2Fidalangref_lwarx_lwri_instrs.htm

⁴⁰ *stcxw* (Store Word Conditional Indexed) instruction, available here:

http://publib.boulder.ibm.com/infocenter/pseries/v5r3/index.jsp?topic=/com.ibm.aix.aixassem/doc/alan_gref/stcxw.htm

⁴¹ *eieio* (Enforce In-Order Execution of I/O) instruction, available here:

http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.aixassem/doc/alangref/idalangref_eieio_instrs.htm

*PowerPC storage model and AIX programming*⁴² describes where synchronization mechanisms must be used to ensure that the code adheres to the Power Architecture. Although this documentation covers how to write compliant code, it does not cover the performance aspect of using the mechanisms.

Unless the code is hand-tuned assembler code, take advantage of the locking services that are provided by the operating system because they are tuned and provide the necessary synchronization mechanisms. *Power Instruction Set Architecture Version 2.07*⁴³ provides assembler programming examples for sharing storage. For more information, see Appendix B, “Performance tooling and empirical performance analysis” on page 199.

Information about this topic, from the OS perspective, is available here:

- ▶ 4.2.6, “Decimal floating point (DFP)” on page 83 (*AIX*)
- ▶ 6.2.6, “Decimal floating point (DFP)” on page 120 (*Linux*)

2.2.10 Fixed-point load and store quadword instructions

The Power architecture provides load and store instructions that operate on quadwords (16-bytes) of storage. The Load Quadword (**lq**) instruction loads an even-odd pair of general purpose registers from the storage that is addressed by the effective address specified by the instruction. The store quadword (**stq**) instruction stores the contents of an even-odd pair of general purpose registers into the storage that is addressed by the effective address specified by the instruction.

Information about this topic, from the processor perspective, is available here:

- ▶ 2.2.10, “Fixed-point load and store quadword instructions” on page 46 (*processor*)

2.2.11 Instruction fusion

POWER8 instruction fusion combines information from two adjacent instructions into one instruction, such that it executes faster than the non-fused instruction. Two forms of fusion are supported for loads with immediate fields that are larger than the allotted 16 bits provided by the Power architecture. This is typically accomplished by an **addi s** instruction to compute the address followed by a load from that address. The two forms of fusion are “Table of content fusion” on page 46 and “Vector load fusion” on page 47.

Capitalizing on instruction fusion

The instruction fusion capabilities of POWER8 are a feature of the processor and do not require special options for the compilers to use them. However, for best performance, **-qtune=pwr8** (XL family) or **-mtune=power8** (GCC) are advised for best use of this feature.

For hand-tuned assembly language code, ensure that the appropriate pattern of code is used and that the two instructions to be fused are adjacent.

Table of content fusion

An example of table of content fusion is:

ADDIS RT, RA, SI

LD RT, RA, DS (eligible instructions are **LD**, **LBZ**, **LHZ**, **LWZ**)

⁴² *PowerPC storage model and AIX programming: What AIX programmers need to know about how their software accesses shared storage*, Michael Lyons, et al, available here:

<http://www.ibm.com/developerworks/systems/articles/powerpc.html>

⁴³ *Power ISA Version 2.07*, available at <https://www.power.org/documentation/power-isa-version-2-07/>

Where the **RT** of the **ADDIS** is the same as **RA** of the **LD** instruction. POWER8 will internally fuse this into a single instruction

Vector load fusion

An example of vector load fusion is:

```
addi RT,0,SI
```

```
lvx VRT, RA, RB (eligible instructions are lxd2x, lxvw4x, lxvdsx, lvebx, lvehx, lvevx, lvx, lxsdx)
```

Where **RT** of **ADDI** is the same as **RB** of the **LVX** instruction and **RA** cannot be zero. POWER8 will internally fuse this into a single instruction

2.2.12 Event-based branches (or user-level fast interrupts)

The event-based branch facility is a hardware facility that generates event-based exceptions when a certain event criteria is met. As an example, this facility allows application programs to enable hardware to change the EA of the next instruction to be executed when certain events occur to an EA specified by the program.

Information about this topic, from the OS perspective, is available here:

- 6.2.7, “Event-based branches” on page 123 (*Linux*)

2.2.13 Power management and system performance

The POWER8 processor has power saving and performance enhancing features that can be used to lower overall energy usage, while yielding higher performance when needed. The following modes can be enabled and modified in order to use these features.

Dynamic Power Saver: Favor Performance

This mode is intended to provide the best performance. If the processor is being used even moderately, the frequency will be raised to the maximum frequency possible to provide the best performance. If the processors are very lightly used, the frequency will be lowered to the minimum frequency, which is potentially far below the nominal shipped frequency, to save energy. Note that the top frequency achieved is based on system type and is affected by environmental conditions. Also note that when running at the maximum frequency, significantly more energy is being consumed, which means this mode can potentially cause an increase in overall energy consumption.

Dynamic Power Saver: Favor Power

This mode is intended to provide the best performance per watt consumed. The processor frequency is adjusted based on the processor utilization to maintain the workload throughput without using more energy than required to do so. At very high processor utilization levels, the frequency will be raised above nominal, just as in the favor performance mode above. Likewise, at very low processor utilization levels, the frequency will be lowered to the minimum frequency. The frequency ranges are the same for the two Dynamic Power Saver modes, but the algorithm that determines which frequency to set is different.

Dynamic Power Saver: Tunable Parameters

The modes just discussed (“Dynamic Power Saver: Favor Performance” and “Dynamic Power Saver: Favor Power”) are tuned to provide both energy savings and performance increases. However, there may be situations where only top performance is of concern, or, conversely,

where peak power consumption is an issue. The tunable parameters can be used to modify the setting of the processor frequency in these modes to meet these various objectives. Note that modifying these parameters should be done only by advanced users. We suggest that, if there are issues that need to be addressed by the Tunable Parameters, IBM should be directly involved in the parameter value selection.

Idle Power Saver

This mode is intended to save the maximum amount of energy when the system is nearly completely idle. When the processors are found to be nearly idle, the frequency of all processors is lowered to the minimum. Additionally, workloads are dispatched onto a smaller number of processor cores so that the other processor cores can be put into a low energy usage state. When processor utilization increases, the process is reversed: The processor frequency is raised back up to nominal, and the workloads are spread out once again over all of the processor cores. There is no performance boosting aspect in this mode, but entering or exiting this mode may affect overall performance. The delay times and utilization levels for entering and exiting this mode can be adjusted to allow for more or less aggressive energy savings.

The controls for all modes listed above are available on the Advanced System Management Interface and are described in more detail in a white paper available at <http://public.dhe.ibm.com/common/ssi/ecm/en/pow03039usen/POW03039USEN.PDF>. Additionally, the appendix of this white paper includes links to other papers that detail the performance benefits and impacts of using these controls.

2.3 Related publications

The publications that are listed in this section are considered suitable for a more detailed discussion of the topics that are covered in this chapter:

- ▶ *AIX dscr_ctl API sample code*, found at this website:
<https://www.power.org/documentation/performance-guide-for-hpc-applications-on-ibm-power-755-system/> (registration required)
- ▶ *AIX Version 7.1 Release Notes*, found at this website:
<http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.ntl/RELNOTES/GI11-9815-00.htm>
See the section, The **dscrctl** command.
- ▶ *Application configuration for large pages*, found at this website:
http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.prftungd/doc/prftungd/config_apps_large_pages.htm
- ▶ *False Sharing*, found at this website:
<http://msdn.microsoft.com/en-us/magazine/cc872851.aspx>
- ▶ *sync (Synchronize) or dcs (Data Cache Synchronize) instruction*, including information about **sync** and **lwsync** (lightweight sync), found at this website:
http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.assem/doc/a_lang_ref/idalangref_sync_dcs_instrs.htm
- ▶ *The Performance of Runtime Data Cache Prefetching in a Dynamic Optimization System*, found at this website:
<http://www.microarch.org/micro36/html/pdf/1u-PerformanceRuntimeData.pdf>

- ▶ *POWER6 Decimal Floating Point (DFP)*, found at this website:
<https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/Power+Systems/page/POWER6+Decimal+Floating+Point+%28DFP%29>
- ▶ *Power ISA Transactional Memory*, found at this website:
<https://www.power.org/documentation/power-isa-transactional-memory/>
- ▶ *Power Architecture ISA 2.06 Stride N prefetch Engines to boost Application's performance*, found at this website:
<https://www.power.org/documentation/whitepaper-on-stride-n-prefetch-feature-of-isa-2-06/> (registration required)
- ▶ *Power ISA Version 2.07*, found at this website:
<https://www.power.org/documentation/power-isa-version-2-07/>

See the following sections:

- Section 3.1: Program Priority Registers
 - Section 3.2: “or” Instruction
 - Section 4.3.4: Program Priority Register
 - Section 4.4.3: OR Instruction
 - Section 5.3.4: Program Priority Register
 - Section 5.4.2: OR Instruction
 - Book I – 4 Floating Point Facility
 - Book I – 5 Decimal Floating Point
 - Book I – 6 Vector Facility
 - Book I – 7 Vector-Scalar Floating Point Operations (VSX)
 - Book I – Chapter 5 Decimal Floating-Point.
 - Book II – 4.2 Data Stream Control Register
 - Book II – 4.3.2 Data Cache Instructions
 - Book II – 4.4 Synchronization Instructions
 - Book II – A.2 Load and Reserve Mnemonics
 - Book II – A.3 Synchronize Mnemonics
 - Book II – Appendix B. Programming Examples for Sharing Storage
 - Book III – 5.7 Storage Addressing
- ▶ *PowerPC storage model and AIX programming: What AIX programmers need to know about how their software accesses shared storage*, found at this website:
<http://www.ibm.com/developerworks/systems/articles/powerpc.html>

See the following sections:

- Power Instruction Set Architecture
 - Section 4.4.3 Memory Barrier Instructions – Synchronize
- ▶ *Product documentation for XL C/C++ for AIX, V12.1 (PDF format)*, found at this website:
<http://www.ibm.com/support/docview.wss?uid=swg27024811>
 - ▶ *Simple performance lock analysis tool (splat)*, found at this website:
http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.prftools/doc/prftools/idprftools_splat.htm
 - ▶ *What makes Apple's PowerPC memcpy so fast?*, found at this website:
<http://stackoverflow.com/questions/1990343/what-makes-apples-powerpc-memcpy-so-fast>
 - ▶ *What programmers need to know about hardware prefetching?*, found at this website:
<http://www.futurechips.org/chip-design-for-all/prefetching.html>



The POWER Hypervisor

This chapter introduces the POWER8 Hypervisor and describes some of the technical details for this product. It covers the following topics:

- ▶ 3.1, “Introduction to the POWER8 Hypervisor” on page 52
- ▶ 3.2, “POWER8 virtualization” on page 53
- ▶ 3.3, “Related publications” on page 61

3.1 Introduction to the POWER8 Hypervisor

Power Virtualization was introduced in POWER5 systems, so there are many reference materials that are available that cover all three resources (CPU, memory, and I/O), virtualization, capacity planning, and virtualization management. Some of these documents are shown in the reference section at the end of this section, which focuses on POWER8 virtualization, usually. As for any workload deployments, capacity planning, selecting the correct set of technologies, and appropriate tuning are critical to deploying high performing workloads. However, in deploying workloads in virtualized environments, there are more aspects to consider, such as consolidation ratio, workload resource usage patterns, and the suitability of a workload to run in a shared resource environment (or latency requirements).

The first step in the virtualization deployment process is to understand if the performance of a workload in a shared resource environment meets customer requirements. If the workload requires consistent performance with stringent latency requirements, then such workloads must be deployed on a dedicated partition rather than on a shared LPAR. The exceptions are where a shared processor pool is not heavily over committed and over-utilized; such workloads can meet stringent requirements in a shared LPAR configuration also.

It is a preferred practice to understand the resource usage of all workloads that are planned for consolidation on a single system, especially when you plan to use a shared resource model, such as shared LPARs, IBM Active Memory™ Sharing, and VIO server technologies. The next step is to use a capacity planning tool that takes virtualization impacts into consideration, such as the IBM Workload Estimator, to estimate capacity for each partition.

One of the goals of virtualization is maximizing usage. This usage can be achieved by consolidating workloads that peak at different times (that is, in a non-overlapping manner, so each workload (or partition) does not have to be sized for peak usage but rather for average usage). At the same time, each workload can grow to consume free resources from the shared pool that belong to other partitions on the system. This situation allows the packing of more partitions (workloads) on a single system, producing a higher consolidation ratio or higher density on the deployed system. A higher consolidation ratio is a key metric to achieve in the data center, as it helps to reduce the total cost of ownership (TCO).

Let us look at a list of key attributes that require consideration when deploying workloads on a shared resource model (virtualization):

- ▶ Levels of variation between average and peak usage of workloads:
 - A large difference between average and peak usage
 - A small difference between average and peak usage
- ▶ Workloads and their peak duration, frequency, and estimate when they potentially peak:
Select workloads that peak at different times (non-overlapping).
- ▶ Workload Service Level Agreement SLA requirements (latency requirements and their tolerance levels).
- ▶ Ratio of active to inactive (mostly idle) partitions on a system.
- ▶ Provisioning and de-provisioning frequency.
- ▶ IBM PowerVM has a richer set of technology options than virtualization on other platforms. It supports dedicated, shared, and a mix of dedicated and shared resource models for each of the system resources, such as processor cores, memory, and I/O:
 - Shared LPAR: Capped versus uncapped.
 - Shared LPAR: Resources over-commit levels to meet the peak usage (the ratio of virtual processors to physical processor entitled capacity).

- Shared LPAR: Weight selection to assign a level of priority to get uncapped capacity (excess cycles to address the peak usage).
- Shared LPAR: Multiple shared pools to address software licensing costs, which prevents a set of partitions from exceeding its capacity consumption.
- Active Memory Sharing: The size of a shared pool is based on active workload memory consumption:
 - Inactive workload memory is used for active workloads, which reduces the memory capacity of the pool.
 - The Active Memory De-duplication option can reduce memory capacity further.
 - AIX file system cache memory is loaned to address memory demands that lead to memory savings.
 - Workload load variation changes active memory consumption, which leads to opportunity for sharing.
- Active Memory Sharing: A shared pool size determines the levels of memory over-commit. Starts without over-commit and is based on workload consumption that reduces the pool.
- Active Memory Expansion: AIX working set memory is compressed.
- Active Memory Sharing and Active Memory Expansion can be deployed on the same workload.
- Active Memory Sharing: VIO server sizing is critical for CPU and memory.
- Virtual Ethernet: An inter-partition communication VLANs option that is used for higher network performance.
- Shared Ethernet versus host Ethernet.
- Virtual disk I/O: Virtual small computer system interface (vSCSI), N_Port ID Virtualization (NPIV), file-backed storage, and storage pool.
- Dynamic resource movement (DLPAR) to adopt to growth.

3.2 POWER8 virtualization

PowerVM hypervisor and the AIX operating system (AIX V6.1 TL 7, AIX V7.1 TL 1 and later versions) on POWER8 implement enhanced affinity in a number of areas to achieve optimized performance for workloads that are running in a virtualized shared processor logical partition (SPLPAR) environment. By using the preferred practices that are described in this guide, customers can attain optimum application performance in a shared resource environment. This guide covers preferred practices in the context of IBM POWER8 systems, so this section can be used as an addendum to other PowerVM preferred practice documents.

3.2.1 Virtual processors

A virtual processor is a unit of a virtual processor resource that is allocated to a partition or virtual machine. PowerVM hypervisor can map a whole physical processor core, or it can create a time slice of a physical processor core.

PowerVM hypervisor creates time slices of Micro-Partitioning on physical CPUs by dispatching and undischatching the various virtual processors for the partitions that are running in the shared pool.

If a partition has multiple virtual processors, they might or might not be scheduled to run simultaneously on the physical processor cores.

Partition entitlement is the guaranteed resource available to a partition. A partition that is defined as capped can consume only the processors units that are explicitly assigned as its entitled capacity. An uncapped partition can consume more than its entitlement, but is limited by many factors:

- ▶ Uncapped partitions can exceed their entitlement if there is unused capacity in the shared pool, dedicated partitions that share their physical processor cores while active or inactive, unassigned physical processors, and Capacity on Demand (CoD) utility processors.
- ▶ If the partition is assigned to a virtual shared processor pool, the capacity for all of the partitions in the virtual shared processor pool might be limited.
- ▶ The number of virtual processors in an uncapped partition is throttled depending on how much CPU it can consume. Here are some examples:
 - An uncapped partition with one virtual CPU can consume only one physical processor core of CPU resources under any circumstances.
 - An uncapped partition with four virtual CPUs can consume only four physical processor cores of CPU.
- ▶ Virtual processors can be added or removed from a partition using HMC actions.

Sizing and configuring virtual processors

The number of virtual processors in each LPAR in the system ought not to *exceed* the number of cores available in the system (central electronic complex (CEC)/framework). Or, if the partition is defined to run in a specific virtual shared processor pool, the number of virtual processors ought not to exceed the maximum that is defined for the specific virtual shared processor pool. Having more virtual processors that are configured than can be running at a single point in time does not provide any additional performance benefit and can actually cause more context switches of the virtual processors, which reduces performance.

If there are sustained periods during which there is sufficient demand for all the shared processing resources in the system or a virtual shared processor pool, it is prudent to configure the number of virtual processors to match the capacity of the system or virtual shared processor pool.

A single virtual processor can consume a whole physical core under two conditions:

- ▶ SPLPAR has an entitlement of 1.0 or more processors.
- ▶ The partition is uncapped and there is idle capacity in the system.

Therefore, there is no need to configure more than one virtual processor to get one physical core.

For example, a shared pool is configured with 16 physical cores. Four SPLPARs are configured, each with entitlement 4.0 cores. To configure virtual processors, consider the sustained peak demand capacity of the workload. If two of the four SPLPARs were to peak to use 16 cores (the maximum available in the pool), then those two SPLPARs would need 16 virtual CPUs. If the other two SPLPARs peak only up to eight cores, those two SPLPARs would be configured with eight virtual CPUs.

Entitlement versus virtual processors

Entitlement is the capacity that an SPLPAR is ensured to get as its share from the shared pool. Uncapped mode allows a partition to receive excess cycles when there are free (unused) cycles in the system.

Entitlement also determines the number of SPLPARs that can be configured for a shared processor pool. The sum of the entitlement of all the SPLPARs cannot exceed the number of physical cores that are configured in a shared pool.

For example, a shared pool has eight cores and 16 SPLPARs are created, each with 0.1 core entitlement and one virtual CPU. We configured the partitions with 0.1 core entitlement because these partitions are not running that frequently. In this example, the sum of the entitlement of all the 16 SPLPARs comes to 1.6 cores. The rest of the 6.4 cores and any unused cycles from the 1.6 entitlement can be dispatched as uncapped cycles.

At the same time, keeping entitlement low when there is capacity in the shared pool is not always a preferred practice. Unless the partitions are frequently idle, or there is a plan to add more partitions, the preferred practice is that the sum of the entitlement of all the SPLPARs configured is close to the capacity in the shared pool. Entitlement cycles are guaranteed, so when a partition is using its entitlement cycles, the partition is not preempted; however, a partition can be preempted when it is dispatched to use excess cycles. Following this preferred practice allows the hypervisor to optimize the affinity of the partition's memory and processor cores and also reduces unnecessary preemptions of the virtual processors.

Entitlement also has effects on the choice of memory and processors that are assigned by the hypervisor for the partition. The hypervisor uses the entitlement value as a guide to the amount of CPU a partition will consume. If the entitlement is undersized, performance can be adversely affected, for example, if there are 4 cores per processor chip and 2 partitions are consistently consuming about 3.5 processors of CPU capacity. If the partitions are undersized with 4 virtual processors and 2.0 entitlement (that is, entitlement is set below normal usage levels), the hypervisor may allocate both of the partitions on the same processor chip, as the entitlement of 2.0 allows two partitions to fit into a 4-core processor chip. If both partitions consistently consume 3.5 processors worth of capacity, the hypervisor will be forced to dispatch some of the virtual processors on chips that do not contain memory associated with the partitions. If the partitions had been configured with an entitled capacity of 3.5 instead of 2.0, the hypervisor would place each partition on its own processor chip to ensure that there is sufficient processor capacity for each partition. This improves the locality, resulting in better performance.

Matching entitlement of an LPAR close to its average usage for better performance

The aggregate entitlement (minimum or wanted processor) capacity of all LPARs in a system is a factor in the number of LPARs that can be allocated. The minimum entitlement is what is needed to boot the LPARs, but the wanted entitlement is what an LPAR gets if there are enough resources available in the system. The preferred practice for LPAR entitlement is to match the entitlement capacity to average usage and let the peak be addressed by more uncapped capacity.

When to add more virtual processors

When there is sustained need for a shared LPAR to use more resources in the system in uncapped mode, increase the virtual processors.

How to estimate the number of virtual processors per uncapped shared LPAR

The first step is to monitor the usage of each partition and for any partition where the average utilization is about 100%, and then add one virtual processor, that is, use the capacity of the configured virtual processors before you add more. Additional virtual processors run concurrently if there are enough free processor cores available in the shared pool.

If the peak usage is below the 50% mark, then there is no need for more virtual processors. In this case, look at the ratio of virtual processors to configured entitlement and if the ratio is greater than 1, then consider reducing the ratio. If there are too many virtual processors configured, AIX can *fold* those virtual processors so that the workload can run on fewer virtual processors to optimize virtual processor performance.

For example, if an SPLPAR is given a CPU entitlement of 2.0 cores and four virtual processors in an uncapped mode, then the hypervisor can dispatch the virtual processors to four physical cores concurrently if there are free cores available in the system. The SPLPAR uses unused cores and the applications can scale up to four cores. However, if the system does not have free cores, then the hypervisor dispatches four virtual processors on two cores so that the concurrency is limited to two cores. In this situation, each virtual processor is dispatched for a reduced time slice as two cores are shared across four virtual processors. This situation can impact performance, so AIX operating system processor folding support might be able to reduce to number of virtual processors that are dispatched so that only two or three virtual processors are dispatched across the two physical cores.

Virtual processor management: Processor folding

The AIX operating system monitors the usage of each virtual processor and the aggregate usage of a shared processor partition to manage the use of virtual processors that are actively engaged by a partition. This management task is carried out using a threshold value that is used to increase, decrease, or hold steady the number of engaged virtual processors for the partition. The threshold is observable as the `vpm_fold_threshold` output by the `schedo` command.

When the aggregate usage goes below the threshold, AIX starts folding down the virtual CPUs so that fewer virtual CPUs are dispatched. This action has the benefit of virtual CPUs running longer before being preempted, which helps improve performance. If a virtual CPU gets a shorter dispatch time slice, then more workloads are cut into time slices on the processor core, which can cause higher cache misses. If the aggregate usage of an SPLPAR goes above the threshold, AIX starts unfolding virtual CPUs so that additional processor capacity can be given to the SPLPAR. AIX cannot engage more virtual processors than are currently defined for the partition. Virtual processor management dynamically adopts the number of virtual processors to match the load on an SPLPAR. This threshold (`vpm_fold_threshold`) represents the SMT thread usage starting with AIX V6.1 TL6. In versions prior to AIX V6.1 TL6, `vpm_fold_threshold` represents the core utilization. The threshold is processor type specific.

When folding increases the number of virtual processors that are engaged, and there are free cores available in the shared processor pool, then unfolding another virtual processor results in the partition getting another core along with its associated caches. Now the partition can run on two primary threads of two cores, instead of two threads (primary and secondary) on the same core. A workload that is running on two primary threads of two cores can achieve higher performance if there is less sharing of data, than the workload that is running on primary and secondary threads of the same core. The AIX virtual processor management default policy aims at using the primary thread of each virtual processor first; therefore, it unfolds the next virtual processor without using the SMT threads of the first virtual processor. After it unfolds all the virtual processors and consumes the primary thread of all the virtual processors, it starts using the secondary and tertiary threads of the virtual processors.

For further details, see this website:

<http://public.dhe.ibm.com/common/ssi/ecm/en/pow03049usen/POW03049USEN.PDF>

Processor bindings in a shared LPAR

In AIX V6.1 TL 7 and in AIX V7.1 TL 1 and later versions, binding virtual processors is available to an application that is running in a shared LPAR. An application process can be bound to a virtual processor in a shared LPAR. In a shared LPAR, a virtual processor is dispatched by the PowerVM hypervisor. The PowerVM hypervisor maintains three levels of affinity for dispatching, such as core, chip, and node level. By maintaining affinity at the hypervisor level and in AIX, applications can achieve higher level affinity through virtual processor bindings.

3.2.2 Page table sizes for LPARs

The hardware page table of an LPAR is sized based on the maximum memory size of an LPAR and not what is assigned (or wanted) to the LPAR. There are some performance considerations if the maximum size is set higher than the wanted memory:

- ▶ A larger page table tends to help performance of the workload, as the hardware page table can hold more pages. This larger table reduces translation page faults. Therefore, if there is enough memory in the system and you want to improve translation page faults, set your max memory to a higher value than the LPAR wanted memory.
- ▶ On the downside, more memory is used for hardware page table, which not only wastes memory, but also makes the table become sparse, which results in the following situations:
 - A dense page table tends to help with better cache affinity because of reloads.
 - Less memory that is consumed by the hypervisor for the hardware page table means that more memory is made available to the applications.
 - There is less page walk time as page tables are small.

3.2.3 Placing LPAR resources to attain higher memory affinity

POWER8 PowerVM optimizes the allocation of resources for both dedicated and shared partitions as each LPAR is activated. Correct planning of the LPAR configuration enhances the possibility of getting both CPU and memory in the same domain in relation to the topology of a system.

PowerVM hypervisor selects the required processor cores and memory that is configured for an LPAR from the system free resource pool. During this selection process, hypervisor takes the topology of the system into consideration and allocates processor cores and memory where both resources are close. This situation ensures that the workload on an LPAR has lower latency in accessing its memory.

When you install a new system, power on the partitions of highest importance first. By doing so, the partitions have first access to available memory and processing resources. After a partition has been powered on, the server has been IPLed, or the Dynamic Platform Optimizer has been run, the processors and memory assignment have been predetermined by the hypervisor, so that the order of activation is not important. On the HMC, there is an option to activate the current configuration, and when you use this option, there will be no change in the current placement of the partition. Activating with a partition profile may change the current placement of the partition.

Partition powering on: Even though a partition is dependent on a VIOS, it is safe to power on the partition before the VIOS. The partition does not fully power on because of its dependency on the VIOS, but claims its memory and processing resources.

How to determine if an LPAR is contained within a domain

From an AIX LPAR, run **lssrad** to display the number of domains across which an LPAR is spread.

The **lssrad** syntax is as follows:

```
lssrad -av
```

If all the cores and memory are in a single domain, you will receive the following output with only one entry under REF1:

REF1	SRAD	MEM	CPU
0	0	31806.31	0-31
	1	31553.75	32-63

REF1 represents a domain, and domains vary by platform. SRAD always references a chip. However, **lssrad** does not report the actual physical domain or chip location of the partition: it is a relative value whose purpose is to inform if the resources of the partition are within the same domain or chip. The output of this **lssrad** example indicates that the LPAR is allocated with 16 cores from two chips within the same domain. Note that the **lssrad** command output was taken from an SMT4 platform, and, thus, CPU 0-31 actually represents 8 cores.

When all the resources are free (an initial machine state or reboot of the CEC), the PowerVM allocates memory and cores as optimally as possible. At partition boot time, PowerVM is aware of all of the LPAR configurations, so placement of processors and memory are made regardless of the order of activation of the LPARs.

However, after the initial configuration, the setup might not stay static. Numerous operations take place, such as:

- ▶ Reconfiguration of existing LPARs with new profiles
- ▶ Reactivating existing LPARs and replacing them with new LPARs
- ▶ Adding and removing resources to LPARs dynamically (DLPAR operations)

Any of these changes might result in memory fragmentation, causing LPARs to be spread across multiple domains. There are ways to minimize or even eliminate the spread. For the first two operations, the spread can be minimized by releasing the resources that are currently assigned to the deactivated LPARs.

Resources of an LPAR can be released by running the following commands:

- ▶ **chhwres -r mem -m <system_name> -o r -q <num_of_Mbytes> --id <lp_id>**
- ▶ **chhwres -r proc -m <system_name> -o r --procunits <number> --id <lp_id>**

The first command frees the memory and the second command frees cores.

Fragmentation because frequent movement of memory or processor cores between partitions is avoidable with correct planning. DLPAR actions can be done in a controlled way so that the performance impact of resource addition or deletion is minimal. Planning for growth helps alleviate the fragmentation that is caused by DLPAR operations. Knowing the LPARs that must grow or shrink dynamically, and placing them with LPARs that can tolerate nodal crossing latency (less critical LPARs), is one approach to handling the changes of critical LPARs dynamically. In such a configuration, when growth is needed for the critical LPAR, the resources that are assigned to the non-critical LPAR can be reduced so that the critical LPAR can grow. Another method of managing fragmentation is to monitor the affinity score of the system or important partitions and use the Dynamic Platform Optimizer to re-optimize the memory and processor assigned to the partitions.

Affinity groups

PowerVM firmware has support for affinity groups that can be used to group multiple LPARs within the same processor chip, processor socket, or drawer. When using affinity groups, it is important to understand the physical configuration of the processor cores and memory contained within the processor chips, processor sockets, and drawers, such that the size of the affinity group does not exceed the capacity of the desired domain.

For example, if the system has 4 cores and 64 GB of memory per processor chip, and you want to contain the partitions to a single processor chip, ensure that the size of the affinity group does not exceed 4 cores and 64 GB of memory. When calculating memory size of an affinity group and what is available on a chip, the computed value needs to account for the memory used by the hypervisor for I/O space and for objects associated with the partition, such as the hardware page table. As a general rule, the size of the affinity group desired memory should only allocate 90% to 95% of the physical memory contained in a domain. If the affinity group is larger than the desired domain, the hypervisor will not be able to contain the affinity group within a single domain.

This affinity group feature can be used in multiple situations:

- ▶ LPARs that are dependent or related, such as server and client, and application server and database server, can be grouped so they are in the same book.
- ▶ Affinity groups can be created that are large enough such that they force the assignment of LPARs to be in different books. For example, if you have a two-socket system and the total resources (memory and processor cores) assigned to the two groups exceeds the capacity of a single socket, these two groups are forced to be in separate sockets.

If a pair of LPARs is created with the intent of one being a failover to another partition, and one partition fails, the other partition (which is placed in the same node, if both are in the same affinity group) uses all of the resources that were freed up from the failed LPAR.

The following HMC CLI command adds or removes a partition from an affinity group:

```
chsyscfg -r prof -m <system_name> -i name=<profile_name>  
lpar_name=<partition_name>,affinity_group_id=<group_id>
```

group_id is a number 1 - 255 (255 groups can be defined), and **affinity_group_id=none** removes a partition from the group.

When the hypervisor places resources at frame reboot, it first places all the LPARs in group 255, then the LPARs in group 254, and so on. Place the most important partitions regarding affinity in the highest configured group.

PowerVM resource consumption for capacity planning considerations

PowerVM hypervisor consumes a portion of the memory resources in the system. During your planning stage, consider the layout of LPARs. Factors that affect the amount of memory that is consumed are the size of the hardware page tables in the partitions, the number of I/O devices, hypervisor memory mirroring, and other factors. Use the *IBM System Planning Tool* to estimate the amount of memory that is reserved by the hypervisor. This tool is available at the following website:

<http://www.ibm.com/systems/support/tools/systemplanningtool/>

Licensing resources and Capacity Upgrade on Demand (CUoD)

Some Power Systems support capacity upgrade on demand so that customers can license capacity on demand as business needs for compute capacity grows. Therefore, a Power System might not have usage of all of the resources that are installed, which poses a challenge to allocate both cores and memory from a local domain. PowerVM correlates customer configurations and licensed resources to allocated cores and memory from the local domain to each of the LPARs.

For systems with unlicensed memory, the licensing is governed on a quantity of memory basis and not on a physical DIMM basis. Therefore, any installed memory can be used to optimize the affinity of partitions. For systems with unlicensed processors, during a CEC reboot, Dynamic Platform Optimizer (see 3.2.5, “Optimizing Resource Placement: Dynamic Platform Optimizer” on page 61), and some DLPAR requests, the hypervisor is able to readjust which processors are licensed to optimize the affinity of the partitions.

For more information about this topic, see 3.3, “Related publications” on page 61.

3.2.4 Active memory expansion

Active memory expansion (AME) is a capability that is supported on POWER8 servers that employs memory compression technology to expand the effective memory capacity of an LPAR. The operating system identifies the least frequently used memory pages and compresses them. The result is that more memory capacity within the LPAR is available to sustain more load, or the ability to remove memory from the LPAR to be used to deploy more LPARs. The POWER8 processor provides enhanced support of AME with the inclusion of on-chip accelerators onto which the work of compression and decompression is offloaded.

AME is deployed by first using the **amepat** tool to model the projected expansion factor and CPU usage of a workload. This modeling looks at the compressibility of the data, the memory reference patterns, and current CPU usage of the workload. AME can then be enabled for the LPAR by setting the expansion factor. The operating system then reports the physical memory available to applications as *actual memory* times the *expansion factor*. Then, transparently, the operating system locates and compresses cold pages to maintain the appearance of expanded memory.

Applications do not need to change, and they are not aware that AME is active. However, not all applications or workloads have suitable characteristics for AME. Here is a partial list of guidelines for the workload characteristics that can be a good fit for AME:

- ▶ The memory footprint is dominated by application working storage (such as heap, stack, and shared memory).
- ▶ Workload data is compressible.
- ▶ Memory access patterns are concentrated to a subset of the overall memory footprint.
- ▶ Workload performance is acceptable without the use of larger page sizes, such as 64 KB pages. AME disables the usage of large pages and uses only 4 KB pages.
- ▶ The average CPU usage of the workload is below 60%.
- ▶ Users of the application and workload are relatively insensitive to response time increases.

For more information about AME usage, see *Active Memory Expansion: Overview and Usage Guide*, available at this website:

<http://public.dhe.ibm.com/common/ssi/ecm/en/pow03037usen/POW03037USEN.PDF>

3.2.5 Optimizing Resource Placement: Dynamic Platform Optimizer

The Dynamic Platform Optimizer feature automates the manual steps to improve resource placement. For more information, visit the following website and select the Doc-type **Word document** *P7 Virtualization Best Practice*. An update to this document from the POWER8 perspective is planned:

<https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/61ad9cf2-c6a3-4d2c-b779-61ff0266d32a/page/64c8d6ed-6421-47b5-a1a7-d798e53e7d9a/attachments>

Note: This document is intended to address POWER8 processor technology-based PowerVM best practices to attain the best LPAR performance. Use this document in conjunction with other PowerVM documents.

3.2.6 Partition compatibility mode

When partitions are created, the processor compatibility mode can be specified. On POWER8, partitions can run in POWER6, POWER6+, POWER7, POWER8, or default compatibility modes. Different modes support different SMT levels and hardware instructions, based on the hardware model that is chosen.

For example, to migrate to a POWER6 server, the partition must be selected to run in POWER6 mode. In addition to allowing migration, the partition, even on a POWER8 server, will run at most in SMT2 mode, and only instructions that are available on POWER6 can be used by the partition. SMT2 mode is used when POWER6 or POWER6+ is selected, SMT4 mode is used for POWER7, and SMT8 mode is used for POWER8, although the default mode for POWER8 on AIX is SMT4. AIX also supports the `smtctl` command, which can increase or reduce the SMT level of the partition if that is desired. A value of `default` means that the partition runs in whatever mode was available when the partition was activated. The selection of the default will prevent the partition from migrating to earlier generations of POWER processors.

3.3 Related publications

The publications that are listed in this section are considered suitable for a more detailed discussion of the topics that are covered in this chapter:

- ▶ *Active Memory Expansion: Overview and Usage Guide*, found here:
<http://public.dhe.ibm.com/common/ssi/ecm/en/pow03037usen/POW03037USEN.PDF>
- ▶ *IBM PowerVM Active Memory Sharing Performance*, found here:
http://public.dhe.ibm.com/common/ssi/rep_wh/n/POW03017USEN/POW03017USEN.PDF
- ▶ *IBM PowerVM Virtualization Introduction and Configuration*, SG24-7940
- ▶ *IBM PowerVM Virtualization Managing and Monitoring*, SG24-7590
- ▶ *POWER7 Virtualization Best Practice Guide*, found here:
https://www.ibm.com/developerworks/wikis/download/attachments/53871915/P7_virtualization_bestpractice.doc?version=1
- ▶ *PowerVM Migration from Physical to Virtual Storage*, SG24-7825

- ▶ *Virtual I/O (VIO) and virtualization*, found here:
<https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/Power%20Systems/page/Virtual%20I%20and%20virtualization>
- ▶ *Virtualization Best Practice*, found here:
<https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/Power%20Systems/page/Virtualization%20best%20practices>



AIX

This chapter describes the optimization and tuning of POWER8 and other Power processor-based servers running the AIX operating system. It covers the following topics:

- ▶ 4.1, “Introduction” on page 64
- ▶ 4.2, “Using Power features with AIX” on page 64
- ▶ 4.3, “AIX operating system-specific optimizations” on page 86
- ▶ 4.4, “AIX preferred practices” on page 96
- ▶ 4.5, “Related publications” on page 98

4.1 Introduction

AIX is regarded as a good choice for building an IT infrastructure on IBM systems that are designed with Power architecture technology. With its proven scalability, advanced virtualization, security, manageability, and reliability features, it is an enterprise-class OS. In particular, AIX is currently the only operating system that leverages decades of IBM technology innovation designed to provide the highest level of performance and reliability of any UNIX operating system. AIX has demonstrated leadership performance on a variety of system benchmarks.

The following performance benefits of AIX are included:

- ▶ Deep integration for Power architecture:
 - Core design with the Power architecture
- ▶ Autonomic optimization:
 - Single OS image configures itself to support any Power processor
 - Dynamic workload optimization
- ▶ Performs on a wide variety of system configurations:
 - Scales from 0.05 to 256 cores (up to 1024 logical processors)
 - Horizontal (native clustering) and vertical scaling
- ▶ Strong virtualization support for PowerVM virtualization:
 - Tight integration with PowerVM
 - Enabler for virtual I/O (VIO)
- ▶ Full set of integrated performance tools

AIX 6.1 and AIX 7.1 run on and maximize on the capabilities of systems based on POWER8, the latest generation of POWER-based technology systems, while supporting POWER4, POWER5, POWER6, and POWER7 (including POWER 7+) systems.

For more information about this topic, see 4.5, “Related publications” on page 98.

4.2 Using Power features with AIX

Some of the significant features of POWER with POWER7 and POWER8 extensions in an AIX environment are described in this section.

4.2.1 Multi-core and multi-thread

Operating system enablement usage of multi-core and multi-thread technology varies by operating system and release. Table 4-1 shows the maximum processor cores and threads for a (single) logical partition running AIX.

Table 4-1 Multi-thread per core features by single LPAR scaling

Single LPAR scaling	AIX release
32-core/32-thread	5.3/6.1/7.1
64-core/128-thread (SMT2)	5.3/6.1/7.1
64-core/256-thread (SMT4)	6.1(TL4)/7.1
256-core/1024-thread (SMT4) (default) or 128-core/1024-thread (SMT8)	7.1

Information about multi-thread per core features by POWER generation is available in Table 2-1 on page 24.

Using multi-core and multi-thread features is a challenging prospect. In addition to the overview material in this section, the following specific scaling topics are described:

- ▶ Malloc tuning (see 4.3.1, “Malloc” on page 86)
- ▶ Pthread tuning (see 4.3.2, “Pthread tunables” on page 89)

Further information about this topic, from the processor and OS perspectives, is available here:

- ▶ 2.2.1, “Multi-core and multi-thread” on page 23 (*processor*)
- ▶ 5.2.1, “Multi-core and multi-thread” on page 102 (*IBM i*)
- ▶ 6.2.1, “Multi-core and multi-thread” on page 108 (*Linux*)

For more information about this topic, see 4.5, “Related publications” on page 98.

Simultaneous Multithreading (SMT)

Simultaneous Multithreading (SMT) is a feature of the Power architecture and is described in “SMT” on page 25. SMT is supported in AIX as described in *Simultaneous multithreading*.¹

AIX provides options to allow SMT customization. The `smtctl` command allows the SMT feature to be enabled, disabled, or capped (SMT2 versus SMT4 mode on POWER7 and SMT2 or SMT4 modes versus SMT8 on POWER8). The partition-wide tuning option, `smtctl`, changes the SMT mode of all processor cores in the partition. It is built on the AIX dynamic reconfiguration (AIX DR) framework to allow hardware threads (logical processors) to be added and removed in a running partition. Because of this option’s global nature, it is normally set by system administrators. Most AIX systems (commercial) use the default SMT settings enabled (that is, SMT2 mode on POWER5 and POWER6, and SMT4 mode on POWER7 and POWER8).

When SMT is enabled (SMT2, SMT4, or SMT8 mode), the AIX kernel takes advantage of the platform feature to dynamically change SMT modes. These mode switches are done based on partition load (the number of running or waiting to run software threads) to choose the optimal SMT mode for the CPUs in the partition. The mode switching policies optimize overall workload throughput, but do not attempt to optimize individual software threads.

Information about the topic of SMT, from the processor and OS perspectives, is available here:

- ▶ “SMT” on page 25 (*processor*)
- ▶ “SMT” on page 102 (*IBM i*)
- ▶ “Simultaneous multithreading (SMT)” on page 109 (*Linux*)

¹ *Simultaneous multithreading*, available here:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.genprogc/doc/genprogc/smt.htm>

SMT priorities

SMT priorities in the Power hardware are introduced in “SMT priorities” on page 25.

AIX kernel usage of SMT thread priority and effects

The AIX kernel is optimized to take advantage of SMT thread priority by lowering the SMT thread priority in select code paths, such as when spinning in the wait process. When the kernel modifies the SMT thread priority and execution is returned to a process-thread, the kernel sets the SMT thread priority back to Medium or the level that is specified by the process-thread using an AIX system call that modified the SMT thread priority (see “Application programming interface (APIs)” on page 66).

Where to use

SMT thread priority can be used to improve the performance of a workload by lowering the SMT thread priority that is being used on an SMT thread that is running a particular process-thread when:

- ▶ The thread is waiting on a lock.
- ▶ The thread is waiting on an event, such as the completion of an IO event.

Alternatively, process-threads that are performance sensitive can maximize their performance by ensuring that the SMT thread priority level is set to an elevated level.

Application programming interface (APIs)

There are three ways to set the SMT priority when it is running on POWER processors:^{2, 3}

- ▶ Modify the SMT priority directly using the PPR register.⁴
- ▶ Modify the SMT priority through the usage of special no-ops.⁵
- ▶ Using the AIX `thread_set_smt_priority` system call.⁶

For more information about this topic, see Table 2-2 on page 26.

Information about the topic of SMT priorities, from the processor and OS perspectives, is available here:

- ▶ “SMT priorities” on page 25 (*processor*)
- ▶ “SMT priorities” on page 110 (*Linux*)

Affinitization and binding

Affinity performance effects are explained in “The POWER8 processor and affinity performance effects” on page 14. Establishing good affinity is accomplished by understanding the placement of a partition on the underlying cores and memory of a Power system, and then by using operating system facilities to bind application threads to run on specific hardware threads or cores.

Flexible SMT

On POWER7 and POWER7+, there is a correlation between the hardware thread number (0-3) and the hardware resources within the processor. Matching the thread numbers to the number of active threads was required for optimum performance. For example, if only one

² Power ISA Version 2.07, available at <https://www.power.org/documentation/power-isa-version-2-07/>

³ `thread_set_smt_priority` or `thread_read_smt_priority` System Call, available here:

http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.kerneltechref/doc/ktechrfl/thread_set_smt_priority.htm

⁴ Power ISA Version 2.07, found here: <https://www.power.org/documentation/power-isa-version-2-07/>

⁵ Ibid

⁶ `thread_set_smt_priority` or `thread_read_smt_priority` System Call, available here:

http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.kerneltechref/doc/ktechrfl/thread_set_smt_priority.htm

thread was active, it was thread0; if two threads were active, they were thread0 and thread1. On POWER8, the same performance is obtained regardless of which thread is active. The processor balances resources according to the number of active threads. There is no need to match the thread numbers with the number of active tasks. Thus, when using the **bindprocessor** command or API, it is not necessary to bind the job to thread0 for optimal performance.

Affinity APIs

Most applications must be bound to logical processors to get a performance benefit from memory affinity to prevent the AIX dispatcher from moving the application to processor cores in different Multi-chip Modules (MCMs) while the application runs.

The most likely way to obtain a benefit from memory affinity is to limit the application to running only on the processor cores that are contained in a single MCM. You can accomplish this task by running the **bindprocessor** command and the **bindprocessor()** function. It can also be done with the resource set affinity commands (**rset**) and service applications. Often, affinity is provided as an administrator option that can be optionally enabled on large systems.

When the application requires more processor cores than contained in a single MCM, the performance benefit through memory affinity depends on the memory allocation and access patterns of the various threads in the application. Applications with threads that individually allocate and reference unique data areas might see improved performance.

Processor affinity (bindprocessor)

Processor affinity is the probability of dispatching of a thread to the logical processor that was previously running it. If a thread is interrupted and later redispached to the same logical processor, the processor's cache might still contain lines that belong to the thread. If the thread is dispatched to a different logical processor, it probably experiences a series of cache misses until its cache working set is retrieved from RAM or the other logical processor's cache. If a dispatchable thread must wait until the logical processor that it was previously running on is available, the thread might experience an even longer delay.

The highest possible degree of processor affinity is to bind a thread to a specific logical processor. Binding means that the thread is dispatched to that logical processor only, regardless of the availability of other logical processors.

The **bindprocessor** command and the **bindprocessor()** subroutine bind the thread (or threads) of a specified process to a particular logical processor. Explicit binding is inherited through **fork()** and **exec()** system calls. The **bindprocessor** command requires the process identifier of the process whose threads are to be bound or unbound, and the bind CPU identifier of the logical processor to be used.

CPU binding is useful for CPU-intensive applications; however, it can sometimes be counter productive for I/O-intensive applications.

RSETS

Every process and kernel thread can have an RSET attached to it. The CPUs on which a thread can be dispatched are controlled by a hierarchy of resource sets. RSETs are mandatory bindings and are honored by the AIX kernel always. Also, RSETs can affect dynamic reconfiguration (DR) activities.

Resource sets

These resource sets are as follows:

Thread effective RSET	Created by ra_attachrset() . Must be a subset (improper or proper) of "Other RSETs" on page 68.
------------------------------	--

Thread partition RSET	Used by WLM. Partition RSETS allow WLM to limit where a thread can run.
Process effective RSET	Created by <code>ra_attachrset()</code> , <code>ra_exec()</code> , and <code>ra_fork()</code> . Must be a subset (improper or proper) of the process partition RSET.
Process partition RSET	Used by WLM to limit where processes in a WLM class are allowed to run. Can also be created by root users using the <code>rs_setpartition()</code> service.

Other RSETs

Another type of RSET is the exclusive RSET. Exclusive use processor resource sets (XRSETs) allow an installation to limit the usage of the processors in XRSETs; they are used only by work that is attached to those XRSETs. They can be created by running the `mkrset` command in the 'sysxrset' namespace.

RSET data types and operations

The public shipped header file `rset.h` contains declarations for the public RSET data types and function prototypes.

An RSET is an opaque data type. Applications allocate an RSET by calling `rs_alloc()`. Applications receive a *handle* to the RSET. The RSET handle (datatype `rsethandle_t` in `sys/rset.h`) is then used in RSET APIs to manipulate or attach the RSET.

Summary of RSET commands

Here is a summary of the RSET commands:

- ▶ **lsrset**: Displays RSETS stored in the system registry or RSETS attached to a process. Here is an example:

lsrset -av	Displays all RSETS in the system registry.
lsrset -p 28026	Displays the effective RSET attached to PID 28026.
- ▶ **mkrset**: Makes a named RSET containing specific CPU and memory pools and place the RSET in the system registry. For example, `mkrset -c 6 10-12 test/lotsofcpus` creates an RSET named `test/lotsofcpus` that contains the specified CPUs.
- ▶ **rmrset**: Removes an RSET from the system registry. Here is an example:


```
rmrset test/lotsofcpus
```
- ▶ **attachrset**: Attaches an RSET to a specified PID. The RSET can either be in the system registry or CPUs or mempools that are specified in the command. Here is an example:

attachrset test/lotsofcpus 28026	Attaches an RSET in a register to a process.
attachrset -c 4-8 28026	Attaches an RSET with CPUs 4 - 8 to a process as an effective RSET.
attachrset -P -c 4-8 28026	Attaches an RSET with CPUs 4 - 8 to process as a partition rset.
- ▶ **detachrset**: Detaches an RSET from a specified PID. Here is an example:

detachrset 28026	Detaches an effective RSET from a PID.
detachrset -P 20828	Detaches a partition RSET from a PID.
- ▶ **execrset**: Runs a specific program or command with a specified RSET. Here is an example:


```
execrset sys/node.04.00000 -e test
```

 Runs a program `test` with an effective RSET from the system registry.

execrset -c 0-1 -e test2	Runs program test2 with an effective RSET that contains logical CPU IDs 0 and 1.
execrset -P -c 0-1 -e test3	Runs program test3 with a partition RSET that contains logical CPU IDs 0 and 1.

RSET manipulation and information services

This list contains only user space APIs. There are also similar kernel extension APIs. For example, **krs_alloc()** is the kernel extension equivalent to **rs_alloc()**.

rs_alloc()	Allocates and initializes an RSET and returns an RSET handle to a caller.
rs_free()	Frees a resource set. The input is an RSET handle.
rs_init()	Initializes a previously allocated RSET. The initialization options are the same as for rs_alloc() .
rs_op()	Performs one of a set of operations against one or two RSETS.
rs_getinfo()	Get information about an RSET.
rs_getrad()	Get resource allocation domain information from an input RSET.
rs_numrads()	Returns the number of system resource allocation domains at the specified system detail level that have available or online resources.
rs_getpartition()	Gets a process's partition RSET.
rs_setpartition()	Sets a process's partition RSET.
rs_discardname()	
rs_getnameattr()	
rs_getnamedrset()	
rs_setnameattr()	
rs_registername()	These are services that are used to manage the RSET system registry. There are services to create, obtain, and delete RSETs in the registry.

Attachment services

Here are the RSET attachment services:

ra_attachrset()	A service to attach a work component to an RSET. The service uses the rstype_t and rsid_t parameters to identify the work component to attach to the input RSET (specified by an rsethandle_t).
ra_detachrset()	Detaches an RSET from the work unit that is specified by the rstype_t/rsid_t parameters.
ra_exec()	Runs a program that is attached to a specific work component. The service uses rstype_t and rsid_t to specify the work component. However, the only supported rstype_t is R_RSET . All of the various versions of exec() are supported.
ra_fork()	Forks a process that is attached to a specific work component. The service uses rstype_t and rsid_t to specify the work component. However, the only supported rstype_t is R_RSET .

ra_get_attachinfo()	The ra_attachrset() also allows RSETs to be attached to ranges of memory in a file or in a shared memory segment.
ra_free_attachinfo()	This service frees the memory that was allocated for the attachment information which was returned by ra_get_attachinfo() .
ra_getrset()	Retrieves the RSET attachment to a process or thread. The return code indicates where the returned RSET is attached.
ra_mmap() and ra_mmapv()	Maps a file or memory region into a process and attaches it to the resource set that is specified by the rstype_t and rsid_t parameters. A memory allocation policy similar to ra_attachrset() allows a caller to specify how memory is preferentially allocated when the area is accessed.
ra_shmget() and ra_shmgetv()	Gets a shared memory segment with an attachment to a resource set. The RSET is specified by the rstype_t and rsid_t parameters. A memory allocation policy similar to ra_attachrset() allows a caller to specify how memory is preferentially allocated when the area is accessed.

AIX Enhanced Affinity (Scheduler Resource Allocation Domain, or SRAD)

AIX Enhanced Affinity is a collection of AIX internal system changes and API extensions to improve performance on IBM POWER7 Systems™. Enhanced Affinity improves performance by increasing CPU and memory locality on POWER7 Systems. Enhanced Affinity extends the AIX existing memory affinity support. AIX V6.1 technology level 6100-05 contains AIX Enhanced Affinity support.

Enhanced Affinity status is determined during system boot and remains unchanged for the life of the system. A reboot is required to change the Enhanced Affinity status. In AIX V6.1.0 technology level 6100-05, Enhanced Affinity is enabled by default on POWER7 machines. Enhanced Affinity is available only on POWER7 machines. Enhanced Affinity is disabled by default on POWER6 and earlier machines. A **vmo** command tunable (**enhanced_memory_affinity**) is available to disable Enhanced Affinity support on POWER7 machines.

The following are two concepts that are related to Enhanced Affinity:

- ▶ **SRAD:** SRAD is the collection of logical CPUs and physical memory resources that are close from a hardware affinity perspective. An AIX system (partition) can consist of one or more SRADs. An SRAD represents the same collection of system resources as an existing MCM. A specific SRAD in a partition is identified by a number. It is an **sradit_t** data type and is often referred to as an SRADID.
- ▶ **SRADID:** The numeric identifier of a specific SRAD. It is a short integer data type. An SRADID value is the index of the resource allocation domain at the R_SRADSDL system detail level in the system's resource set topology.

Power Systems before POWER7 Systems provided only system topology information to dedicated CPU logical partitions. This setup limited the usefulness of RSET attachments for CPU and memory locality purposes to dedicated CPU partitions. POWER7 Systems provide system topology information for shared CPU logical partitions (SPLPAR).

The **1ssrad** command can be used to display the processor and memory resources at the SRAD and REF1 levels (where REF1 is the next higher level affinity domain).

You can use the AIX Enhanced Affinity services to attach SRADs to threads and memory ranges so that the application preferentially identifies the logical CPUs or physical memory to use to run the application. AIX continues to support RSET attachments to identify resources for an application.

RSET versus SRADs

When you compare RSET with SRADIDs, observe these considerations:

1. SRADIDs can be attached to threads, shared memory segments, memory map regions, and process memory subranges. SRADIDs might not be attached at the process level (R_PROCESS). SRADIDs might not be attached to files (R_FILDES).
2. SRADID attachments are considered advisory. There are no mandatory SRADID attachments. AIX might ignore advisory SRADID attachments.
3. Process and thread RSET attachments continue to be mandatory. The process and thread resource set hierarchy continues to be enforced. Memory RSET attachments (shared memory, file, and process subrange) continue to be advisory. This situation is unchanged from previous affinity support.

API support

SRADIDs can be attached to threads and memory by using the following functions:

- ▶ `ra_attach()` (new)
- ▶ `ra_fork()`
- ▶ `ra_exec()`
- ▶ `ra_mmap()` and `ra_mmapv()`
- ▶ `ra_shmget()` and `ra_shmgetv()`

SRADIDs can be detached from thread and memory by using the `sra_detach()` function (new).

Information about the topic of affinization and binding, from the processor and OS perspectives, is available here:

- ▶ “Affinitization and binding to hardware threads” on page 26 (*processor*)
- ▶ “Affinitization and binding” on page 111 (*Linux*)

Hybrid thread and core

AIX provides facilities to customize SMT characteristics of CPUs running within a partition. The features require some partition-wide CPU configuration options, so their use is limited to specific workloads.

Hybrid thread features

AIX provides some basic features that allow more control in SMT mode. With these features, specific software threads can be bound to hardware threads assigned to ST mode CPUs. This configuration allows for an asymmetric SMT configuration, where some CPUs are in high SMT mode, and others have SMT mode disabled. This configuration allows critical software threads within a workload to receive an ST performance boost, and allows the remaining threads to benefit from SMT mode. Here are some typical reasons to take advantage of this hybrid mode:

- ▶ Asymmetric workload, where the performance of one thread serializes an entire workload. For example, one master thread dispatches work to many subordinate threads.
- ▶ Software threads that are critical to a system administrator.

The ability to create hybrid SMT configurations is limited under current AIX releases and does require administrator or privileged configuration changes. CPUs that provide ST mode hardware threads must be placed into XRSETs. XRSETs contain logical CPUs that are segregated from the general kernel dispatching pool. Software threads must be explicitly bound to CPUs in an XRSET. The only way to create an XRSET is by running the **mkrset** command. All of the hardware threads for logical CPUs must be contained in the XRSET created RSET. To accomplish this task, run the following commands:

lsrset -av	Displays the RSET topology. The system CPU topology is broken down into a hierarchy that has the form <code>sys/node.XX.YYYYY</code> . The largest XX value is the CPU (core) level. This command provides logical processor groups by core.
mkrset -c 4-7 sysxrset/set1	Creates an XRSET <code>sysxrset/set1</code> containing logical CPUs 4 - 7.

An XRSET alone can be used to ensure that only specific work uses a CPU set. There is also the ability to restrict work execution to primary threads in an XRSET. This ability is known as an STRSET. STRSETs allow software threads to use ST execution mode independently of the load on the other CPUs in the system. Work can be placed onto STRSETs by running the following commands:

execrset -S	This command allows external programs to start and be bound to an exclusive RSET.
ra_attach(R_STRSET)	This API allows a thread to be bound to an STRSET.

Information about this topic, from the processor and OS perspectives, is available here:

- “Hybrid thread and core” on page 27 (*processor*)
- “Hybrid thread and core” on page 112 (*Linux*)

For more information about this topic, see 4.5, “Related publications” on page 98.

AIX folding

Folding is a key AIX feature on shared processor LPARs that can improve both system and partition performance. Folding is needed for supporting a large number of partitions in a system. It is an integrated feature, requiring both hardware and PowerVM support. The AIX component that manages folding is the Virtual Processor Manager (VPM).

The basic concept of folding is to compress work to a smaller number of cores, based on CPU utilization, by folding the remaining cores. The unused cores are folded by VPM, and PowerVM does not schedule them for dispatch in the partition unless the operating system requests that the cores be unfolded (or woken up), for example, when the workload changes or when a timer interrupt needs to be fired on that core.

As an example, an LPAR might have 24 virtual cores (processors) assigned, but is only consuming a total of three physical processors across all of these virtual cores. Folding compresses (moves) all work to a smaller number of cores (three cores plus some extra cores to handle spikes in workload), allowing PowerVM to allocate the unused cores for use elsewhere on the system.

Folding generally improves LPAR and system performance by reducing context switching of cores between partitions across a system, thus reducing context switching of software threads across multiple cores in a LPAR. It improves overall affinity at both the LPAR and system levels.

VPM runs once per second and computes how many cores will be kept unfolded based on the overall CPU utilization of the LPAR. On POWER8 systems, the folding algorithm has been enhanced to include the average load (or the average number of runnable software threads) as a factor in the computation.

Folding can be enabled and disabled using the **schedo** command to adjust the value of the **vpm_fold_policy** tunable. To respond faster to spikes in workloads or on partitions with a high interrupt load, a second tunable, **vpm_xvcpus**, can also be used to increase the number of spare, unfolded CPUs. This can improve response time for workloads with steep utilization spikes, or on partitions with a high interrupt load, although this can result in higher core usage.

AIX V6.1 TL8 and AIX V7.1 TL2 introduced a new scaled throughput-based folding algorithm that can be enabled using the **schedo** command to adjust the value of the **vpm_throughput_mode** tunable. The default folding algorithm favors single-threaded performance and overall LPAR throughput over core utilization. The new scaled throughput algorithm can favor reduced core utilization and higher core throughput, instead of overall LPAR throughput. The new algorithm applies both load and utilization data to make folding decisions. It can switch unfolded cores to SMT2, SMT4, or SMT8 modes when the workload increases, rather than unfolding more cores. This is depicted in Figure 4-1.

The degree of SMT mode, SMT2 or SMT4 (or SMT8 for POWER8 systems), to favor reduced core utilization can be controlled by assigning the appropriate value to the **vpm_throughput_mode** tunable (2 for SMT2 mode, 4 for SMT4 mode, and 8 for SMT8 mode). When the **vpm_throughput_mode** tunable is set to a value of 1, the folding algorithm behaves like the legacy folding algorithm and favors single-threaded (ST mode) performance. However, unlike the legacy algorithm, which uses only utilization data, the new algorithm employs both load and utilization data to make folding decisions.

The default value of the **vpm_throughput_mode** tunable is 0 (zero) on POWER8 systems, just like on POWER7 and earlier systems (the legacy folding algorithm will continue to be applicable).

If the **vpm_throughput_mode** is set to a value of 1 or greater, then the **vpm_throughput_core_threshold** tunable can also be set to specify the number of cores that must be unfolded before the **vpm_throughput_mode** parameter is honored. One scheme that balances between performance and core utilization when enabling higher SMT modes is to set **vpm_throughput_core_threshold** to the integer value of the entitled capacity.

The scaled throughput algorithm can reduce overall core utilization at the frame level for certain workloads (Figure 4-1).

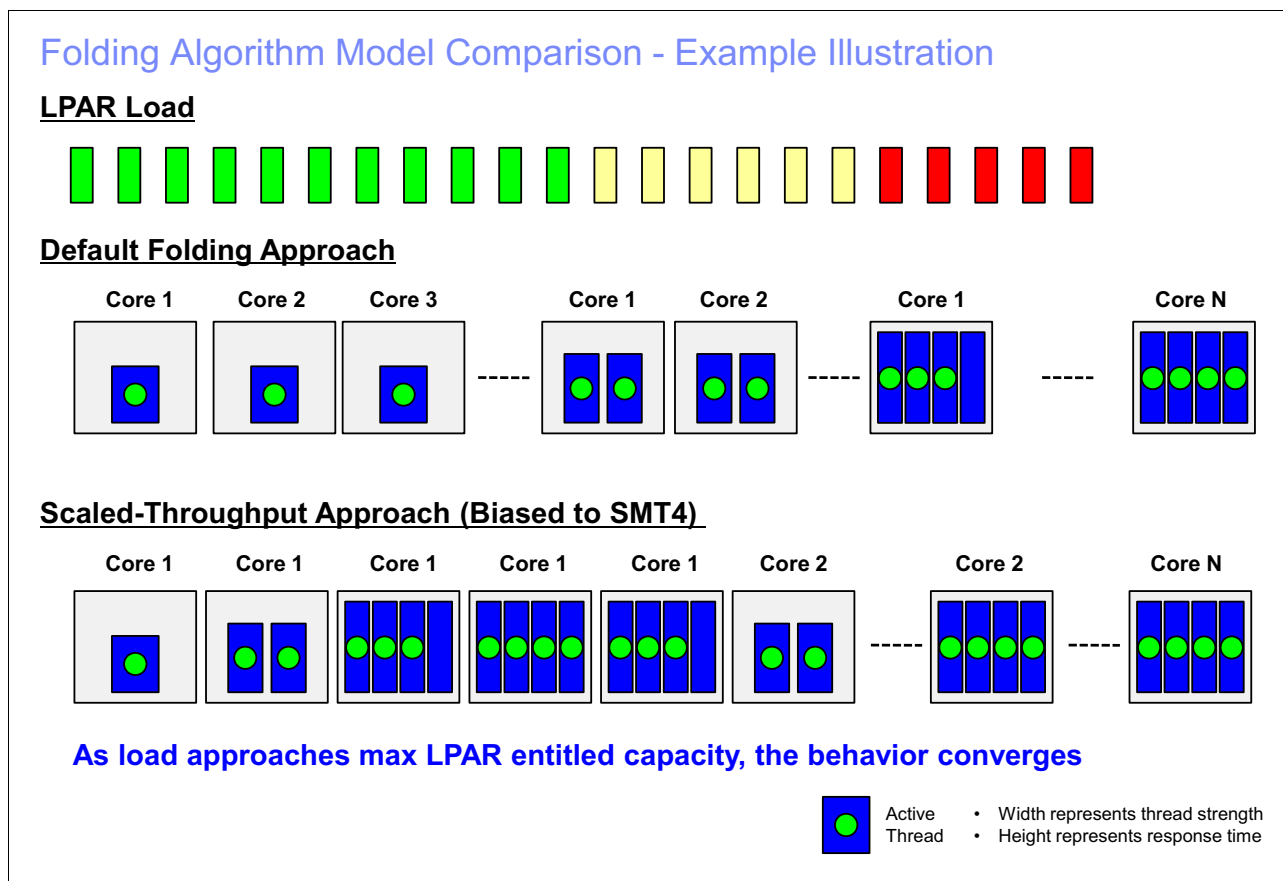


Figure 4-1 Folding algorithm model comparison

4.2.2 Multipage size support on AIX

AIX supports up to four different page sizes (see Table 4-2), but the actual page sizes that are supported by a particular system vary, based on processor chip type. The **pagesize -a** command on AIX determines all of the page sizes that are supported by AIX on a particular system.

Because the 64 KB page size is easy to use, and because it is expected that many applications perform better when they use the 64 KB page size rather than the 4 KB page size, AIX has rich support for the 64 KB page size. No system configuration changes are necessary to enable a system to use the 64 KB page size. On systems that support the 64 KB page size, the AIX kernel automatically configures the system to use it. Table 4-2 and Table 4-3 list the page size specifications for Power Systems.

Table 4-2 Page size support for Power HW and AIX configuration support⁷

Page size	Required hardware	Requires user configuration	Restricted
4 KB	ALL	No	No
64 KB	POWER5+ or later	No	No
16 MB	POWER4 or later	Yes	Yes
16 GB	POWER5+ or later	Yes	Yes

Table 4-3 Supported segment page sizes on AIX⁸

Segment base page size	Supported page sizes	Minimum required hardware
4 KB	4 KB/64 KB	POWER6
64 KB	64 KB	POWER5+
16 MB	16 MB	POWER4
16 GB	16 GB	POWER5+

Page sizes are an attribute of an individual segment. Earlier POWER processors only supported a single page size per segment. The system administrator or user had to choose the optimal page size for a specific application based on its memory footprint. POWER 5+ introduced the concept of mixed or multiple page sizes within a single segment: 4 KB and 64 KB. POWER7 and later processors support mixed page segment sizes of 4 KB, 64 KB and 16 MB.

Starting with version 6.1, AIX takes advantage of this new hardware capability on POWER6 and later processors to combine the conservative memory usage aspects of the 4 KB page size in sparsely referenced memory regions, with the performance benefits of the 64 KB page size in densely referenced memory regions. AIX V6.1 takes advantage of this automatically, without user intervention, although it is disabled in segments that have an explicit page size selected by the user. This AIX feature is referred to as dynamic Variable Page Size Support (VPSS). Some applications might prefer to use a larger page size, even when a 64 KB region is not fully referenced. The page size promotion aggressiveness factor (PSPA) can be used to reduce the memory-referenced requirement, at which point a group of 4 KB pages is promoted to a 64 KB page size. The `vmo` command on AIX allows configuration of the VMM tunable parameters. The PSPA can be set for the whole system by using the `vmm_default_pspa` `vmo` tunable, or for a specific process by using the `vm_pattr` system call.⁹

In addition to 4 KB and 64 KB page sizes, AIX supports 16 MB pages, also called *large pages*, and 16 GB pages, also called *huge pages*. These page sizes are intended for use only in high-performance environments, and AIX by default does not automatically configure a system to use these page sizes.

Use the `vmo` tunables `lgpg_regions` and `lgpg_size` to configure the number of 16 MB large pages on a system.

The following example allocates 1 GB of 16 MB large pages:

```
vmo -r -o lgpg_regions=64 -o lgpg_size=16777216
```

⁷ Multiple page size support, available here:

http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.prftungd/doc/prftungd/multiple_page_size_support.htm

⁸ Ibid

⁹ Ibid

To use large pages, non-root users must have the CAP_BYPASS_RAC_VMM capability in AIX enabled. The system administrator can add this capability by running **chuser**:

```
chuser capabilities=CAP_BYPASS_RAC_VMM,CAP_PROPAGATE <user_id>
```

Huge pages must be configured using the Hardware Management Console (HMC). To do so, complete the following steps:

1. On the managed system, click **Properties** → **Memory** → **Advanced Options** → **Show Details** to change the number of 16 GB pages.
2. Assign 16 GB huge pages to a partition by changing the partition profile.

The **vmo** tunable **vmm_mpsize_support** can be used to limit multiple page size support. The default value of 1 supports all four page sizes, but the tunable can be set to other values to configure which page sizes are to be supported.

Application support to use multisize pages on AIX¹⁰

As described in *Power Instruction Set Architecture Version 2.07*,¹¹ you can specify page sizes to use for four regions of a 32-bit or 64-bit process address space.

These page sizes can be configured with an environment variable or with settings in an application XCOFF binary with the **ldedit** or **ld** commands, as shown in Table 4-4.

Table 4-4 Page sizes for four regions of a 32-bit or 64-bit process address space

Region	ld or ldedit option	LDR_CNTRL environment variable	Description
Data	bdatapsize	DATAPSIZE	Initialized data, bss, and heap
Stack	bstacksize	STACKSIZE	Initial thread stack
Text	btextpsize	TEXTPSIZE	Main executable text
Shared memory	None	SHMPsize	Shared memory that is allocated by the process

You can specify a different page size to use for each of the four regions of a process address space. Only the 4 KB and 64 KB page sizes are supported for all four memory regions. The 16 MB page size is supported only for the process data, process text, and process shared memory regions. The 16 GB page size is supported only for a process shared memory region.

You can set the preferred page sizes for an application in the XCOFF/XCOFF64 binary file by running the **ldedit** or **ld** commands.

The **ld** or **cc** commands can be used to set these page size options when you are linking an executable command:

- **ld -o mpsize.out -btextpsize:4K -bstacksize:64K sub1.o sub2.o**
- **cc -o mpsize.out -btextpsize:4K -bstacksize:64K sub1.o sub2.o**

The **ldedit** command can be used to set these page size options in an existing executable command:

```
ldedit -btextpsize=4K -bdatapsize=64K -bstacksize=64K mpsize.out
```

¹⁰ Ibid

¹¹ *Power ISA Version 2.07*, available at <https://www.power.org/documentation/power-isa-version-2-07/>

We can set the preferred page sizes of a process with the **LDR_CNTRL** environment variable. As an example, the following command causes the `mpsize.out` process to use 4 KB pages for its data, 64 KB pages for its text, 64 KB pages for its stack, and 64 KB pages for its shared memory on supported hardware:

```
LDR_CNTRL=DATAPSIZE=4K@TEXTFSIZE=64K@SHMPSIZE=64K mpsize.out
```

Page size environment variables override any page size settings in an executable XCOFF header. Also, the **DATAPSIZE** environment variable overrides any **LARGE_PAGE_DATA** environment variable setting.

Rather than using the **LDR_CNTRL** environment variable, consider marking specific executable files to use large pages, because this limits the large page usage to the specific application that benefits from large page usage.

Page size and shared memory

To back shared memory segments of an application with large pages, specify the **SHM_LGPAGE** and **SHM_PIN** flags in the `shmget()` function. In addition, set the **vmo v_pinshm** tunable to a value of 1 with, for example, `vmo -r -o v_pinshm=1`. If large pages are unavailable, the 4 KB pages back the shared memory segment.

Support for specifying the page size to use for the shared memory of a process with the **SHMPSIZE** environment variable is available starting in IBM AIX 5L™ Version 5.3 with the 5300-08 Technology Level, or later, and AIX Version 6.1 with the 6100-01 Technology Level, or later.

Monitoring page size that is used by an application

Monitoring page size is accomplished by running the following commands:¹²

- ▶ The **ps** command can be used to monitor the base page sizes that are used for process data, stack, and text.
- ▶ The **vmstat** command has two options available to display memory statistics for a specific page size:
 - The **vmstat -p** command displays global **vmstat** information, along with a breakdown of statistics per page size.
 - The **vmstat -P** command displays per page size statistics.

Information about this topic, from the processor and OS perspectives, is available here:

- ▶ 2.2.2, “Multipage size support: Page sizes (4 KB, 64 KB, 16 MB, and 16 GB)” on page 27 (*processor*)
- ▶ 4.2.2, “Multipage size support on AIX” on page 74
- ▶ 5.2.2, “Multipage size support on IBM i” on page 103
- ▶ 6.2.2, “Multipage size support on Linux” on page 113

For more information about this topic, see 4.5, “Related publications” on page 98.

¹² *Multiple page size support*, available here:

http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.prftungd/doc/prftungd/multiple_page_size_support.htm

4.2.3 Efficient use of cache

Generally, with Power Architecture, unlike some other architectures, users do not need to be too concerned about cache management or optimizing cache usage. This section describes AIX facilities for controlling hardware prefetching through the Data Streams Control Register (DSCR) and is meant for advanced users who understand their workload characteristics and want to experiment with the register settings for improving performance. See 2.2.3, “Efficient use of cache and memory” on page 28 for a discussion that includes a more detailed description of the DSCR register and its settings.

Controlling Data Streams Control Register (DSCR) under AIX

Under AIX, Data Streams Control Register (DSCR) settings can be controlled both by the `dscr_ctl()` programming API and by running the `dscrctl` command.^{13,14}

`dscr_ctl()` API

```
#include <sys/machine.h>
int dscr_ctl(int op, void *buf_p, int size)
```

Where:

- op:** Operation. Possible values are `DSCR_WRITE`, `DSCR_READ`, `DSCR_GET_PROPERTIES`, and `DSCR_SET_DEFAULT`.
- Buf_p:** Pointer to an area of memory where the values are copied from (`DSCR_WRITE`) or copied to (`DSCR_READ` and `DSCR_GET_PROPERTIES`). For `DSCR_WRITE`, `DSCR_READ`, and `DSCR_SET_DEFAULT` operations, `buf_p` must be a pointer to a 64-bit data area (long long *). For `DSCR_GET_PROPERTIES`, `buf_p` must be a pointer to a `struct dscr_properties` (defined in `<sys/machine.h>`).
- Size:** Size in bytes of the area pointed to by `buf_p`.

Function:

The action that is taken depends on the value of the operation parameter that is defined in `<sys/machine.h>`:

- | | |
|----------------------------|--|
| DSCR_WRITE | Stores a new value from the input buffer into the process context and in the DSCR. |
| DSCR_READ | Reads the current value of DSCR and returns it in the output buffer. |
| DSCR_GET_PROPERTIES | Reads the number of hardware streams that are supported by the platform, the platform (firmware) default Prefetch Depth and the Operating System default Prefetch Depth from kernel memory, and returns the values in the output buffer (<code>struct dscr_properties</code> defined in <code><sys/machine.h></code>). |
| DSCR_SET_DEFAULT | Sets a 64-bit DSCR value in a buffer pointed to by <code>buf_p</code> as the operating system default. Returns the old default in the buffer pointed to by <code>buf_p</code> . Requires root authority. The new default value is used by all the processes that do not explicitly set a DSCR value using <code>DSCR_WRITE</code> . The new default is not permanent across reboots. For an operating system default prefetch depth that is permanent across reboots, use the <code>dscrctl</code> command, which adds an entry into the <code>inittab</code> to initialize the system-wide prefetch depth default value upon reboot (for a description of this command, see “The <code>dscrctl</code> command” on page 80). |

¹³ `dscr_ctl` Subroutine, available here:

http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.basetechref/doc/basetrf1/dscr_ctl.htm

¹⁴ `dscrctl` Command, available here:

<http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.cmds/doc/aixcmds2/dscrctl.htm>

Return values are as follows:

0 if it is successful.

-1 if an error is detected. In this case, `errno` is set to indicate the error. Here are some possible values:

EINVAL	Invalid value for DSCR (DSCR_WRITE, DSCR_SET_DEFAULT).
EFAULT	Invalid address that is passed to function.
EPERM	Operation not permitted (DSCR_SET_DEFAULT by non-root user).
ENOTSUP	Data streams that are not supported by platform hardware.

Symbolic values for the following SSE and DPFD fields are defined in `<sys/machine.h>`:

DPFD_DEFAULT	0
DPFD_NONE	1
DPFD_SHALLOWEST	2
DPFD_SHALLOW	3
DPFD_MEDIUM	4
DPFD_DEEP	5
DPFD_DEEPER	6
DPFD_DEEPEST	7
DSCR_SSE	8

Here is a description of the **dscr_properties** structure in `<sys/machine.h>`:

```
struct dscr_properties {
    uintversion;
    uintnumber_of_streams; /* Number of HW streams */
    longlongplatform_default_pd; /* PFW default */
    longlongos_default_pd; /* AIX default */
    longlong dscr_res[5]; /* Reservd for future use */
};
```

Here is an example of this structure:

```
#include <sys/machine.h>
int rc;
long long dscr = DSCR_SSE | DPFD_DEEPER;
rc = dscr_ctl(DSCR_WRITE, &dscr);
...
```

A new process inherits the DSCR from its parent during a *fork*. This value is reset to the system default during *exec*.

When a thread is dispatched (starts running on a CPU), the value of the DSCR for the owning process is written in the DCSR. You do not need to save the value of the register in the process context when the thread is *undispatched* because the system call writes the new value both in the process context and in the DCSR.

When a thread runs **dscr_ctl** to change the prefetch depth for the process, the new value is written into the AIX process context and the DCSR register of the thread that is running the system call. If another thread in the process is concurrently running on another CPU, it starts using the new DSCR value only after the new value is reloaded from the process context area after either an interrupt or a redispatch. This action can take as much as 10 ms (a clock tick).

The dsccrctl command

The system administrator can use this command to read the current settings for the hardware streams mechanism and set a system wide value for the DSCR. The DSCR is privileged. It can be read or written only by the operating system.

To query the characteristics of the hardware streams on the system, run this command:

```
dsccrctl -q
```

Here is an example of this command:

```
# dsccrctl -q
Current DSCR settings:
    number_of_streams = 16
    platform_default_pd = 0x5 (DPFD_DEEP)
    os_default_pd = 0xd (DSCR_SSE | DPFD_DEEP)
```

To set the operating system default prefetch depth on the system temporarily (that is, for the current session) or permanently (that is, after each reboot), run the following command:

```
dsccrctl [-n] [-b] -s <dscr_value>
```

The **dscr_value** is treated as a decimal number unless it starts with 0x, in which case it is treated as hexadecimal.

To cancel a permanent setting of the operating system default prefetch depth at boot time, run the following command:

```
dsccrctl -c
```

Applications that have predictable data access patterns, such as numerical applications that process arrays of data in a sequential manner, benefit from aggressive data prefetching. These applications must run with the default operating system prefetch depth, or whichever settings are empirically found to be the most beneficial.

Applications that have considerably unpredictable data access patterns, such as some transactional applications, can be negatively affected by aggressive data prefetching. The data that is prefetched is unlikely to be needed, and the prefetching uses system bandwidth and might displace useful data from the caches. Some WebSphere Application Server and DB2 workloads have this characteristic. Performance can be improved by disabling hardware prefetching in these cases by running the following command:

```
dsccrctl -n -s 1
```

This system (partition) wide disabling is only appropriate if it is expected to benefit all of the applications that are running in the partition. However, the same effect can be achieved on an application-specific basis by using the programming API.

Information about the efficient use of cache, from the processor and OS perspectives, is available here:

- ▶ 2.2.3, “Efficient use of cache and memory” on page 28 (*processor*)
- ▶ 6.2.3, “Efficient use of cache” on page 113 (*Linux*)

For more information about this topic, see 4.5, “Related publications” on page 98.

4.2.4 Transactional memory (TM)

Transactional memory (TM) is a POWER8 shared-memory synchronization construct that allows process-threads to perform storage operations that appear to be atomic to other process-threads and applications. One of the main uses of TM is that it speeds up the lock-based programs through the speculative execution of lock-based, critical sections, and it does so without first acquiring a lock. This allows applications that have not been carefully tuned for performance to take advantage of the benefits of fine-grain locking. The transactional programming model also provides productivity gains when developing lock-based, shared memory programs.

Although POWER8 supports TM, you need to explicitly check for support of TM before using the facility, because the processor might be running in a compatibility mode, or the operating system or hypervisor might not support the use of TM. In AIX, the preferred API that determines if TM is supported is the **getsystemcfg()** system call. A new **SC_TM_VER** system variable setting is provided that reports whether TM is supported. A new **__power_tm()** macro is provided that allows the caller to determine if TM is supported. Refer to `/usr/include/sys/systemcfg.h`.

Software failure handler

Upon transaction failure, hardware re-directs control to the failure handler associated with the outermost transaction. The discussion about “Transaction failure” on page 38 explains this and provides details about how control is passed to the software failure handler and the machine state of the status registers.

The Power Architecture Platform reserves a range of failure codes for the hypervisor, for client operating systems, and for user applications, to indicate a failure reason when issuing a **tabort.** instruction. These codes are as follows:

- ▶ **0x00 – 0x3F** is reserved for use by AIX
- ▶ **0x40 – 0xDF** is free for use by problem state (application) code
- ▶ **0xE0 – 0xFF** is reserved for use by a hypervisor

The failure codes reserved by AIX to indicate the cause of the failure are defined in `/usr/include/sys/machine.h`.

Debugger support

The dbx AIX debugger, located in `/usr/ccs/bin/dbx`, supports machine-level debugging of TM programs. This support includes the ability to disassemble the new TM instructions, and to display the TM SPRs.

Setting a breakpoint inside of a transaction causes the transaction to unconditionally fail whenever the breakpoint is encountered. To determine the cause and location of a failing transaction, the approach is to set a breakpoint on the transaction failure handler, and then to view the TEXASR and TFIAR registers when the breakpoint is encountered.

The TEXASR, TFIAR, and TFHAR registers can be displayed using the **print** subcommand with the **\$texasr**, **\$tfiar**, or **\$tfhar** parameter. The line of code associated with the address found in TFIAR and TFHAR can be displayed using the **list** subcommand. Here is an example:

```
(dbx) list at $tfiar
```

A new **tm_status** subcommand is provided that displays and interprets the contents of the TEXASR register. This is useful in determining the nature of a transaction failure.

Tracing support

The AIX trace facility has been expanded to include a set of trace events for TM operations performed by AIX, including the processing of TM-type facility unavailable interrupts, preemptions that cause transaction failure, and other operations that can cause transaction failure. The trace event identifier **675** can be used as input to the **trace** and **trcrpt** commands to view TM-related trace events.

System call support

When a system call is made while a processor or thread is transactional (and the transaction has not been suspended), the system call will not be invoked by the AIX kernel. The associated transaction will be persistently failed, and the system call handler will return control to the calling code with an error code of **ENOSYS**. When this occurs, the FC field of the TEXASR register will contain the failure code **TM_ILL_SC**, which is defined in `/usr/sys/include/machine.h`.

It is assumed that any operations performed under a suspended transaction (when the application programmer has explicitly suspended the transaction) are intended to be persistent. Any operations performed by a system call made while in suspended state will not be rolled-back in the event that the transaction fails.

The reason that AIX cannot allow system calls to be made while in transactional state, is that any operations (writes or updates, including I/O) performed by AIX underneath a system call cannot be rolled back.

AIX threads library support

The use of TM is not supported for applications utilizing M:N threads. Undefined behavior might be encountered by transactional threads in an environment where more than one thread shares a single kernel thread. Usage of TM by an application that utilizes M:N threads can lead to a persistent transaction failure with the failure code **TM_PTH_PREEMPTED** being set in TEXASR.

Support of context management subroutines

The use of the context management subroutines, such as the **libc** subroutines **getcontext()**, **setcontext()**, **makecontext()**, **swapcontext()**, **setjmp()**, and **longjmp()** are not supported while in transactional or suspended state. Such operations, where non-transactional context is restored while in transactional or suspended state or context, is saved off while in transactional or suspended state, and then restored while in non-transactional state, leads to an inconsistent state and can result in undefined behavior. Under certain circumstances, AIX fails a transaction attempting to call such routines.

Information about the topic of transactional memory, from the processor, OS, and compiler perspectives, is available here:

- ▶ 2.2.4, “Transactional memory (TM)” on page 37 (*processor*)
- ▶ 6.2.4, “Transactional memory (TM)” on page 113 (*Linux*)
- ▶ 7.3.5, “Transactional memory (TM)” on page 149 (*XL and GCC compiler families*)

4.2.5 Vector Scalar eXtension (VSX)

A program can determine whether a system supports the vector extension by reading the `vmx_version` field of the `_system_configuration` structure. If this field is non-zero, then the system processor chips and operating system contain support for the vector extension. A value of 1 means that the processor chips on the system are Vector Multimedia eXtension (VMX) capable, and a value of 2 means that they are both VMX and Vector Scalar eXtension (VSX) capable. Alternately, the `__power_vmx()` and `__power_vsx()` macros provided in `/usr/include/sys/systemcfg.h` can be used to perform these tests.

Vector capability support in AIX

The AIX Application Binary Interface (ABI) is extended to support the addition of vector register state and conventions. AIX supports the AltiVec programming interface specification.

A set of malloc subroutines (`vec_malloc`, `vec_free`, `vec_realloc`, and `vec_calloc`) is provided by AIX that give 16-byte aligned allocations. Vector-enabled compilation, with `_VEC_` implicitly defined by the compiler, result in any calls to older mallocs and callocs being redirected to their vector-safe counterparts, `vec_malloc` and `vec_calloc`. Non-vector code can also be explicitly compiled to pick up these same malloc and calloc redirections by explicitly defining `__AIXVEC`.

The alignment of the default `malloc()`, `realloc()`, and `calloc()` allocations can also be controlled at run time. This task can be done externally to any program by using the **MALLOCALIGN** environment variable, or internally to a program by using the `mallopt()` interface command option.¹⁵

Information about the topic of VSX, from the processor, OS, and compiler perspectives, is available here:

- ▶ 2.2.5, “Vector Scalar eXtension (VSX)” on page 40 (*processor*)
- ▶ 5.2.3, “Vector Scalar eXtension (VSX)” on page 103 (*IBM i*)
- ▶ 6.2.5, “Vector Scalar eXtension (VSX)” on page 120 (*Linux*)
- ▶ 7.3.2, “Compiler support for VSX” on page 145 (*XL and GCC compiler families*)

For more information about this topic, see 4.5, “Related publications” on page 98.

4.2.6 Decimal floating point (DFP)

Decimal (base 10) data is widely used in commercial and financial applications. However, most computer systems have only binary (base two) arithmetic. There are two binary number systems in computers: integer (fixed-point) and floating point. Unfortunately, decimal calculations cannot be directly implemented with binary floating point. For example, the value 0.1 needs an infinitely recurring binary fraction, whereas a decimal number system can represent it exactly, as one tenth. So, using binary floating point cannot ensure that results are the same as those results using decimal arithmetic.

In general, decimal floating point (DFP) operations are emulated with binary fixed-point integers. Decimal numbers are traditionally held in a binary-coded decimal (BCD) format. Although BCD provides sufficient accuracy for decimal calculation, it imposes a heavy cost in performance, because it is usually implemented in software.

IBM Power processor-based systems, starting with POWER6, provide hardware support for DFP arithmetic. Their microprocessor cores include a DFP unit that provides acceleration for the DFP arithmetic. The IBM Power instruction set is expanded: 54 new instructions were added to support the DFP unit architecture. DFP can provide a performance boost for applications that are using BCD calculations.

¹⁵ *AIX vector programming*, available here:

http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.genprogc/doc/genprogc/vector_prog.htm

How to take advantage of DFP unit on POWER

You can take advantage of the DFP unit on POWER with the following features:¹⁶

- Native DFP language support with a compiler

The C draft standard includes the following new data types (these are native data types, as are int, long, float, double, and so on):

<code>_Decimal32</code>	7 decimal digits of accuracy
<code>_Decimal64</code>	16 decimal digits of accuracy
<code>_Decimal128</code>	34 decimal digits of accuracy

Note: The `printf()` function uses new options to print these new data types:

- `_Decimal32` uses `%Hf`
- `_Decimal64` uses `%Df`
- `_Decimal128` uses `%DDf`

- The IBM XL C/C++ Compiler, release 9 or later, includes native DFP language support. Here is a list of compiler options for IBM XL compilers that are related to DFP:

- **-qdfp:** Enables DFP support. This option makes the compiler recognize DFP literal suffixes, and the `_Decimal32`, `_Decimal64`, and `_Decimal128` keywords.
- **-qfloat=dfpemulate:** Instructs the compiler to use calls to library functions to handle DFP computation, regardless of the architecture level. You might experience performance degradation when you use software emulation.
- **-qfloat=nodfpemulate** (the default when the **-qarch** flag specifies POWER6 or POWER7 or POWER8): Instructs the compiler to use DFP hardware instructions.
- **-D__STDC_WANT_DEC_FP__:** Enables the referencing of DFP-defined symbols.

For hardware supported DFP, with **-qarch=pwr6**, **-qarch=pwr7**, or **-qarch=pwr8**, use the following command:

```
cc -qdfp
```

For software emulation of DFP (on earlier processor chips), use the following command:

```
cc -qdfp -qfloat=dfpemulate
```

- The GCC compilers for Power Systems also include native DFP language support.

The following list describes GCC compiler options that are related to DFP:

- **-mhard-dfp** (the default when **-mcpu=power6** or **-mcpu=power7** is specified): Instructs the compiler to directly take advantage of DFP hardware instructions for decimal arithmetic.
- **-mno-hard-dfp:** Instructs the compiler to use calls to library functions to handle DFP computation, regardless of the architecture level. If your application is dynamically linked to the **libdfp** variant and running on POWER6 or POWER7 processors, then the run time automatically binds to the **libdfp** variant implemented with hardware DFP instructions. Otherwise, the software DFP library is used. You might experience performance degradation when you use software emulation.
- **-D__STDC_WANT_DEC_FP__:** Enables the reference of DFP defined symbols.

¹⁶ How to compile DFPAL?, available here: <http://speleotrove.com/decimal/dfpal/compile.html>

- Decimal Floating Point Abstraction Layer (DFPAL), which is a no additional cost, downloadable library from IBM.¹⁷

Many applications that are using BCD today use a library to perform math functions. Changing to a native data type can be hard work, after which you might have an issue with one code set for AIX on POWER6, POWER7, or POWER8 and one for other platforms that do not support native DFP. The solution to this problem is DFPAL, which is an alternative to the native support. DFPAL contains a header file to include in your code and the DFPAL library.

The header file is downloadable from General Decimal Arithmetic at <http://speleotrove.com/decimal/> (search for “DFPAL”).

Download the complete source code, and compile it on your system.

If you have hardware support for DFP, use the library to access the functions.

If you do not have hardware support (or want to compare the hardware and software emulation), you can force the use of software emulation by setting a shell variable before you run your application by running the following command:

```
export DFPAL_EXE_MODE=DNSW
```

Determining if your applications are using DFP

There are two AIX commands that are used for monitoring:

- **hpmstat** (for monitoring the whole system)
- **hpmcount** (for monitoring a single program)

The **PM_DFU_FIN** (DFU instruction finish) field in the output of the **hpmstat** and **hpmcount** commands verifies that the DFP operations finished.

The **-E PM_MRK_DFU_FIN** option in the **tprof** command uses the AIX trace subsystem, which tells you which functions are using DFP and how often.

Information about this topic, from the processor and OS perspectives, is available here:

- 2.2.6, “Decimal floating point” on page 42 (*processor*)
- 5.2.4, “Decimal floating point” on page 103 (*IBM i*)
- 6.2.6, “Decimal floating point (DFP)” on page 120 (*Linux*)

For more information about this topic, see 4.5, “Related publications” on page 98.

4.2.7 On-chip encryption accelerator

When the AIX operating system runs on POWER7+ or POWER8 processors, it transparently uses on-chip encryption accelerators. For each of the uses that are described in this section, there are no application visible changes or awareness required.

AIX encrypted file system (EFS)

Integrated with the AIX Journaled File System (JFS2) is the ability to create an encrypted file system (EFS) where all data at rest in the file system is encrypted. When AIX EFS runs on POWER7+ or POWER8, it uses the encryption accelerators, which can show up to a 40% advantage in file system I/O-intensive operations. Applications do not need to be aware of this situation, but application and workload deployments might be able to take advantage of higher levels of security by using AIX EFS for sensitive data.

¹⁷ Ibid

AIX Internet Protocol Security (IPSec)

When IPSec is enabled on AIX running on POWER7+ or POWER8, AIX transparently uses the on-chip encryption accelerators for all data in transit. The advantage that is provided by the accelerators is more pronounced when jumbo frames (a maximum transmission unit (MTU) of 9000 bytes) are used. Applications do not need to be aware of this situation, but application and workload deployments might be able to take advantage of higher levels of security by enabling IPSec.

AIX /dev/random (random number generation)

AIX capitalizes on the on-chip random number generator on POWER7+ and POWER8. Applications that use the AIX special files `/dev/random` or `/dev/urandom` transparently get the advantages of stronger hardware-based random numbers. If an application is making high frequency usage of random number generation, there can also be a performance advantage.

AIX PKCS11 Library

On POWER7+ and POWER8 systems, the AIX operating system PKCS11 library transparently uses the on-chip encryption accelerators. For an application using the PKCS11 APIs, no change or awareness by the application is required. The AIX library interfaces dynamically decides, based on the algorithm and data size, when to use the accelerators. Because of the cost of setup and programming of the on-chip accelerators, the advantage is limited to operations on large blocks of data (tens to hundreds of kilobytes).

Information about this topic, from the processor perspective, is available here:

- 2.2.8, “On-chip accelerators” on page 44 (*processor*)

4.3 AIX operating system-specific optimizations

Here we describe optimization methods specific to AIX.

4.3.1 Malloc

Every application needs a fast, scalable, and memory efficient allocator. However, each application’s memory request patterns are different. It is difficult to provide one common allocator or tunable that can satisfy the needs of all applications. AIX provides different memory allocators and suboptions within the allocator, so that a system administrator or developer can choose more suitable settings for their application. This chapter explains the available choices and when to choose them.

Memory allocators

AIX provides three different allocators, and each of them uses a different memory management algorithm and data structures. These allocators work independently, so the application developer must choose one of them by exporting the **MALLOCTYPE** environment variable. The allocators are as follows:

- Default allocator:

The default allocator is selected when the **MALLOCTYPE** environment variable is unset. This setting maintains a consistent performance, even in a worst case scenario, but might not be as memory efficient as a Watson allocator. This allocator is ideal for 32-bit applications, which do not make frequent calls to `malloc()`.

► Watson allocator:

This allocator is selected when `MALLOCTYPE=watson` is set. This allocator is designed for 64-bit applications. It is memory efficient, scalable, and provides good performance. This allocator has a built-in bucket component for allocation requests up to 512 bytes. Table 4-5 provides the mapping for the allocation requests to bucket size.

Table 4-5 Mapping for allocation requests to bucket size

Request size	Bucket size	Request size	Bucket size	Request size	Bucket size	Request size	Bucket size
1 - 4		33-40	40	129-144	144	257-288	288
5 - 8		41 - 48	48	145 - 160	160	289 - 320	320
9 - 12	12	49 - 56	56	161 - 176	176	321 - 352	352
13 - 16	16	57 - 64	64	177 - 192	192	353 - 384	384
17 - 20	20	65 - 80	80	193 - 208	208	385 - 416	416
21 - 24	24	81 - 96	96	209 - 224	224	417 - 448	448
25 - 28	28	97 - 112	112	224 - 240	240	449 - 480	480
29 - 32	32	113 - 128	128	241 - 256	256	481 - 512	512

This allocator is ideal for 64-bit memory-intensive applications.

► Malloc 3.1 allocator:

This allocator is selected when `MALLOCTYPE=3.1` is set. This is a bucket allocator that divides the heap into 28 hash buckets, each with a size of $2^{\text{pow}(x+4)}$, where x stands for bucket index. This allocator provides the best performance at the cost of memory. In most cases, this algorithm can use as much as twice the amount of memory that is actually requested by the application. In addition, an extra page is required for buckets larger than 4096 bytes because objects of a page in size or larger are page-aligned. Interestingly, some earlier customer applications still use this allocator, as it is more tolerant for application memory overwrite bugs.

Memory allocator suboptions

There are many suboptions available that can be selected by exporting the `MALLOCOPTIONS` environment variable. This chapter covers a few of the suboptions that are more relevant to performance tuning. For a complete list of options, see *System memory allocation using the malloc subsystem*, available here:

http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.genprogc/doc/genprogc/sys_mem_alloc.htm

► Multiheap:

By default, the malloc subsystem uses a single heap, which causes lock contention for internal locks that are used by malloc in case of multi-threaded applications. By enabling this option, you can configure the number of parallel heaps to be used by allocators. You can set the multiheap by exporting `MALLOCOPTIONS=multiheap[:n]`, where n can vary between 1- 32 and 32 is the default if n is not specified.

Use this option for multi-threaded applications, as it can improve performance.

► Buckets:

This suboption is similar to the built-in bucket allocator of the Watson allocator. However, with this option, you can have fine-grained control over the number of buckets, number of

blocks per bucket, and the size of each bucket. This option also provides a way to view the usage statistics of each bucket, which be used to refine the bucket settings.

In case the application has many requests of the same size, then the bucket allocator can be configured to preallocate the required size by correctly specifying the bucket options. The block size can go beyond 512 bytes, compared to the Watson allocator or malloc pool options.

You can enable the buckets allocator by exporting `MALLOCOPTIONS=buckets`. Complete details about the buckets options for fine-grained control are available¹⁸. Enabling the buckets allocator turns off the built-in bucket component if the Watson allocator is used.

► **malloc pools:**

This option enables a high performance front end to malloc subsystem for managing storage objects smaller than 513 bytes. This suboption is similar to the built-in bucket allocator of the Watson allocator. However, this suboptions maintains the bucket for each thread, providing lock-free allocation and deallocation for blocks smaller than 513 bytes. This suboption improves the performance for multi-threaded applications, as the time spent on locking is avoided for blocks smaller than 513 bytes.

The pool option makes small memory block allocations fast (no locking) and memory efficient (no header on each allocation object). The pool malloc both speeds up single threaded applications and improves the scalability of multi-threaded applications.

► **malloc disclaim:**

By enabling this option, **free()** automatically disclaims memory. This suboption is useful for reducing the paging space requirement. This option can be set by exporting `MALLOCOPTIONS=disclaim`.

Use cases

Here are some uses cases that you can use to set up your environment:

1. For a 32-bit single-threaded application, use the default allocator.
2. For a 64-bit application, use the Watson allocator.
3. Multi-threaded applications use the **multiheap** option. Set the number of heaps proportional to the number of threads in the application.
4. For single-threaded or multi-threaded applications that make frequent allocation and deallocation of memory blocks smaller than 513, use the **malloc pool** option.
5. For a memory usage pattern of the application that shows high usage of memory blocks of the same size (or sizes that can fall to common block size in bucket option) and sizes greater than 512 bytes, use the **configure malloc bucket** option.
6. For older applications that require high performance and do not have memory fragmentation issues, use **malloc 3.1**.
7. Ideally, the Watson allocator, along with the **multiheap** and **malloc pool** options, is good for most multi-threaded applications. The pool front end is fast and scalable for small allocations, and, in conjunction with **multiheap**, ensures scalability for larger and less frequent allocations.
8. If you notice high memory usage in the application process even after you run **free()**, the **disclaim** option can help.

For more information about this topic, see 4.5, “Related publications” on page 98.

¹⁸ *System memory allocation using the malloc subsystem*, available here:
http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.genprogc/doc/genprogc/sys_mem_alloc.htm

4.3.2 Pthread tunables

The AIX pthread library can be customized with a set of environment variables. Specific variables that improve scaling and CPU usage are listed here. A full description is provided in the following documentation settings:

► **AIXTHREAD_SCOPE={P|S}**

The **P** option signifies a process-wide contention scope (M:N), and the **S** option signifies a system-wide contention scope (1:1). Use system scope (1:1) for AIX. Although process scope (M:N) continues to be supported, it is no longer being enhanced in AIX.

► **SPINLOOPTIME=*n***

The **SPINLOOPTIME** variable controls the number of times the system tries to get a busy mutex or spin lock without taking a secondary action, such as calling the kernel to yield the process. This control is intended for MP systems, where it is hoped that the lock that is held by another actively running pthread is released. The parameter works only within libpthreads (user threads). If locks are usually available within a short period, you might want to increase the spin time by setting this environment variable. The number of times to try a busy lock before yielding to another pthread is *n*. The default is 40 and *n* must be a positive value.

► **YIELDLOOPTIME=*n***

The **YIELDLOOPTIME** variable controls the number of times that the system yields the logical processor when it tries to acquire a busy mutex or spin lock before it goes to sleep on the lock. The logical processor is yielded to another kernel thread, assuming that there is another executable thread with sufficient priority. This variable is effective in complex applications, where multiple locks are in use. The number of times to yield the logical processor before blocking on a busy lock is *n*. The default is 0 and *n* must be a positive value.

For more information about this topic, see 4.5, “Related publications” on page 98.

4.3.3 pollset

AIX 5L V5.3 introduced the pollset APIs. Pollsets are an AIX replacement for UNIX **select()** and **poll()**. **Pollset**, **select()**, and **poll()** all allow an application to efficiently query the status of file descriptors. This action is typically done to allow a single application to multiplex I/O across many file descriptors. Pollset APIs can be more efficient when the number of file descriptors that are queried becomes large.

Efficient I/O event polling through the pollset interface on AIX contains a pollset summary and outlines the most advantageous use of Java. To see this topic, go to this website:

<http://www.ibm.com/developerworks/aix/library/au-pollset/index.html>

For more information about this topic, see 4.5, “Related publications” on page 98.

4.3.4 File system performance benefits

AIX JFS2 is the default file system for 64-bit kernel environments. Applications can capitalize on the features of JFS2 for better performance.

4.3.5 Direct I/O

The AIX read-ahead and write-behind JFS2 feature might not be suitable for applications that perform large sized I/O operations, as the cache hit ratio is low. In those cases, an application developer must evaluate Direct I/O for I/O intensive applications.

Programs that are good candidates for direct I/O are typically CPU-limited and perform much disk I/O. Technical applications that have large sequential I/Os are good candidates. Applications that benefit from striping are also good candidates.

The direct I/O access method bypasses the file cache and transfers data directly from disk into the user space buffer, as opposed to using the normal cache policy of placing pages in kernel memory.

At the user level, file systems can be mounted using the **dio** option on the **mount** command.

At the programming level, applications enable direct I/O access to a file by passing the **O_DIRECT** flag to the open subroutine. This flag is defined in the `fcntl.h` file. Applications must be compiled with **_ALL_SOURCE** enabled to see the definition of **O_DIRECT**.

For more information, see *Working with file I/O*, available here:

http://publib.boulder.ibm.com/infocenter/aix/v6r1/index.jsp?topic=%2Fcom.ibm.aix.gprog%2Fdoc%2Fgenprog%2Fworking_file_io.htm

4.3.6 Concurrent I/O (CIO)

An AIX JFS2 inode lock imposes write serialization at the file level. Serializing write accesses prevents data inconsistency because of overlapping writes. Serializing reads regarding writes ensures that the application does not read stale data.

However, some applications can choose to implement their own data serialization, usually at a finer level of granularity than the file. Therefore, they do not need the file system to implement this serialization for them. The inode lock hinders performance in such cases, by unnecessarily serializing non-competing data accesses. For such applications, AIX offers the concurrent I/O (CIO) option. Under CIO, multiple threads can simultaneously perform reads and writes on a shared file. For applications that do not enforce serialization for accesses to shared files, we suggest not using CIO, as it can result in data corruption because of competing accesses.

Enhanced JFS supports concurrent file access to files. Similar to direct I/O, this access method bypasses the file cache and transfers data directly from disk into the user space buffer.

CIO can be specified for a file either by running **mount -o cio** or by using the **open()** system call (by using **O_CIO** as the **OFlag** parameter).

4.3.7 Asynchronous I/O

If an application does a synchronous I/O operation, it must wait for the I/O to complete. In contrast, asynchronous I/O operations run in the background and do not block user applications, which improves performance, because I/O operations and applications processing can run simultaneously. Many applications, such as databases and file servers, take advantage of the ability to overlap processing and I/O.

Applications can use the `aio_read()`, `aio_write()`, or `lio_listio()` subroutines (or their 64-bit counterparts) to perform asynchronous disk I/O. Control returns to the application from the subroutine when the request is queued. The application can then continue processing while the disk operation is being performed.

4.3.8 I/O completion ports

A limitation of the AIO interface that is used in a threaded environment is that `aio_nwait()` collects completed I/O requests for *all* threads in the same process. One thread collects completed I/O requests that are submitted by another thread.

Another limitation is that multiple threads cannot invoke the collection routines (such as `aio_nwait()`) at the same time. If one thread issues `aio_nwait()` when another thread is calling it, the second `aio_nwait()` returns EBUSY. This limitation can affect I/O performance when many I/Os must run at the same time and a single thread cannot run fast enough to collect all the completed I/Os.

On AIX, using I/O completion ports with AIO requests provides the capability for an application to capture the results of various AIO operations on a per-thread basis in a multi-threaded environment. This functionality provides threads with a method of receiving a completion status for only the AIO requests initiated by the thread.

You can enable IOCP on AIX by running `smitty iocp`. Verify that IOCP is enabled by running the following command:

```
lsdev -Cc iocp
```

The resulting output is shown in the following example:

```
iocp0 Available I/O Completion Ports
```

4.3.9 shmat versus mmap

Memory mapped files provide a mechanism for a process to access files by directly incorporating file data into the process address space. The use of mapped files can reduce I/O data movement because the file data does not have to be copied into process data buffers, as is done by the read and write subroutines. When more than one process maps the same file, its contents are shared among them, providing a low-impact mechanism by which processes can synchronize and communicate.

AIX provides two methods for mapping files and anonymous memory regions. The first set of services, which are known collectively as the *shmat* services, are typically used to create and use shared memory segments from a program. The second set of services, which are known collectively as the *mmap* services, is typically used for mapping files, although it can be used for creating shared memory segments as well.

Both the *mmap* and *shmat* services provide the capability for multiple processes to map the same region of an object so that they share addressability to that object. However, the *mmap* subroutine extends this capability beyond that provided by the *shmat* subroutine by allowing a relatively unlimited number of such mappings to be established. Although this capability increases the number of mappings that are supported per file object or memory segment, it can prove inefficient for applications in which many processes map the same file data into their address space. The *mmap* subroutine provides a unique object address for each process that maps to an object. The software accomplishes this task by providing each process with a unique virtual address, which is known as an *alias*. The *shmat* subroutine allows processes to share the addresses of the mapped objects.

shmat can be used to share memory segments in a way that is similar to how it creates and uses files. An *extended shmat* capability is available for 32-bit applications with their limited address spaces. If you define the **EXTSHM=ON** environment variable, then processes running in that environment can create and attach more than 11 shared memory segments.

Use the shmat services under the following circumstances:

- ▶ When mapping files larger than 256 MB
- ▶ When mapping shared memory regions that must be shared among unrelated processes (no parent-child relationship)
- ▶ When mapping entire files

In general, shmat is more efficient but less flexible.

Use mmap under the following circumstances:

- ▶ Many files are mapped simultaneously.
- ▶ Only a portion of a file must be mapped.
- ▶ Page-level protection must be set on the mapping (allows a 4K boundary).

For more information, see *General Programming Concepts: Writing and Debugging Programs*, available here:

http://publib16.boulder.ibm.com/doc_link/en_US/a_doc_lib/aixprgdd/genprogc/understanding_mem_mapping.htm

For more information about this topic, see 4.5, “Related publications” on page 98.

4.3.10 Large segment tunable aliasing (LSA)

AIX V6.1 TL5 and AIX V7.1 introduce the 1 TB Segment Aliasing. 1 TB segments can improve the performance of 64-bit large memory applications. The optimization is specific to large shared memory (**shmat()** and **mmap()**) regions.

One TB segments are a feature present in POWER5+ and later processors. They can be used to reduce Segment Lookaside Buffer (SLB) misses, and increase the reach of the SLB, reducing the impact of effective-to-virtual address to real translation impact. Applications that are 64-bit and that have large shared memory regions can benefit from incorporating 1 TB segments. This feature is enabled by default on AIX V7.1, but can be enabled using the **vmo** command to adjust the **esid_allocator** tunable.

An overview of 1 TB segment usage can be found in the *IBM AIX Version 7.1 Differences Guide*, SG24-7910.

For more information about this topic, see 4.5, “Related publications” on page 98.

4.3.11 64-bit versus 32-bit ABIs

AIX provides complete support for both 32-bit and 64-bit ABIs. Applications can be developed using either ABI with some performance trade-offs. The 64-bit ABI provides more scaling benefits. With both ABIs, there are performance trade-offs to be considered.

Overview of 64-bit/32-bit ABI

All current POWER processors support a 32-bit and 64-bit execution mode. The 32-bit execution mode is a subset of the 64-bit execution mode. The modes are similar, where the most significant difference is addresses in address generation (effective addresses are truncated to 32 bits) and computation of some fixed-point status registers (carry, overflow, and so on). Although hardware 32-bit/64-bit mode does not affect performance, the 32-bit/64-bit ABIs provided by AIX do have performance implications and tradeoffs.

The 32-bit ABI provides an ILP32 model (32-bit integers, longs, and pointers). The 64-bit ABI provides an LP64 model (32-bit integer and 64-bit longs/pointers). Although current POWER CPUs have 64-bit fixed-point registers, they are treated as 32-bit fixed-point registers by the ABI (the high 32 bits of all fixed-point registers are treated as volatile or undefined by the ABI). The 32-bit ABI preserves only 32-bit fixed-point context across subroutine linkage, non-local goto (**longjmp()**), or signal delivery. 32-bit programs cannot attempt to use 64-bit registers when they run in 32-bit mode (32-bit ABI). In general, other registers (floating point, vector, and status registers) are the same size in both 32-bit/64-bit ABIs.

Starting with AIX V6.1 all supervisor code (kernel, kernel extensions, and device drivers) uses the 64-bit ABI. In general, a unified system call interface is provided to applications that provides efficient system call linkage to both 32-bit and 64-bit applications. Because the AIX V 6.1 kernel is 64-bit, it implies that all systems supported by AIX V 6.1 support the 64-bit ABI. Some older IBM PowerPC CPUs supported on AIX 5L V 5.3 cannot run the 64-bit ABI.

Operating system libraries provide both 32-bit and 64-bit objects, allowing full support for either ABI. Development tools (assembler, linker, and debuggers) support both ABIs.

Trade-offs

The primary motivation to choose the 64-bit ABI is to go beyond the 4 GB directly memory addressability barrier. A second reason is to improve scalability by extending some 32-bit data type limits that are in the 32-bit ABI (`time_t`, `pid_t`, and `offset_t`). Lastly, 64-bit mode provides access to 64-bit fixed-point registers and instructions that can improve the performance of specific fixed-point operations (long long arithmetic and 64-bit memory copies).

The 64-bit ABI does have some performance drawbacks, such as the 64-bit fixed-point registers and the LP64 model grow stack usage and data structures. These items can cause a performance drawback for some applications. Also, 64-bit text is generally larger for most compiles, producing a larger i-cache footprint.

The most significant issue is typically the porting effort (for existing applications), as changing between ILP32 and LP64 normally requires a port. Large memory addressability and scalability are normally the deciding factor when you chose an application execution model.

For more information about this topic, see 4.5, “Related publications” on page 98.

4.3.12 Sleep and wake-up primitives (`thread_wait` and `thread_post`)

AIX provides proprietary `thread_wait()` and `thread_post()` APIs that can be used to optimize thread synchronization and communication in instructions per cycle (IPC). AIX also provides several standard APIs that can be used for thread synchronization and communication. These APIs include `pthread_cond_wait()`, `pthread_cond_signal()`, and `semop()`. Although many applications use these standard APIs, the low-level primitives are available to optimize these operations. `thread_wait()` and `thread_post()` can be used to optimize critical applications services, such as user-mode locking or message passing. They are more efficient than the portable/standard APIs.

Here is more information about the associated subroutines:

► **thread_wait**

The **thread_wait** subroutine allows a thread to wait or block until another thread posts it with the **thread_post** or the **thread_post_many** subroutine or until the time limit specified by the timeout value expires.

If the event for which the thread is waiting and for which it is posted occurs only in the future, the **thread_wait** subroutine can be called with a timeout value of 0 to clear any pending posts by running the following command:

```
thread_wait (timeout)
```

► **thread_post**

The **thread_post** subroutine posts the thread whose thread ID is indicated by the value of the **tid** parameter, of the occurrence of an event. If the posted thread is waiting in **thread_wait**, it is awakened immediately. If it is not waiting in **thread_wait**, the next call to **thread_wait** is not blocked, but returns with success immediately.

Multiple posts to the same thread without an intervening wait by the specified thread counts only as a single post. The posting remains in effect until the indicated thread calls the **thread_wait** subroutine, upon which the posting is cleared.

► **thread_post_many**

The **thread_post_many** subroutine posts one or more threads of the occurrence of the event. The number of threads to be posted is specified by the value of the **nthreads** parameter, and the **tidp** parameter points to an array of thread IDs of threads that must be posted. The subroutine works just like the **thread_post** subroutine, but can be used to post to multiple threads at the same time. A maximum of 512 threads can be posted in one call to the **thread_post_many** subroutine.

For more information about this topic, see 4.5, “Related publications” on page 98.

4.3.13 Shared versus private loads

You can use AIX to share text for libraries and dynamically loaded modules. File permissions can be used to enable and disable sharing of loaded text.

Documentation

AIX provides optimizations that enable sharing of loaded text (libraries and dynamically loaded modules). Sharing text among processes often improves performance because it reduces resource usage (memory and disk space). It also allows unrelated software-threads to share cache space when they run concurrently. Lastly, it can reduce load times when the code is already loaded by a previous program.

Applications can control if private or shared loads are performed to shared text regions. Shared loads require that execute permissions be set for group/other on the text files. As a preferred practice, enable sharing.

For more information about this topic, see 4.5, “Related publications” on page 98.

4.3.14 Workload partitions (WPARs) shared License Program Product (LPP) installs

Starting with AIX V6.1, the WPAR feature gives the system administrator the ability to easily create an isolated AIX operating system that can run services and applications. WPAR provides a secure and isolated environment for enterprise applications in terms of process, signal, and file system space. Any software that is running within the context of a workload partition appears to have its own separate instance of AIX.

The usage of multiple virtual operating systems within a single global operating environment can have multiple advantages. It increases administrative efficiency by reducing the number of AIX instances that must be maintained.

Applications can be installed in a shared environment or a non-shared environment. When an application is installed in a shared environment, it means that it is installed in the global environment and then the application is shared with one or more WPARs. When an application is installed in a non-shared environment, it means that it is installed in the WPAR only. Other WPARs do not have access to that application.

Shared WPAR installation

A shared installation is straightforward because installing software in the global environment is accomplished in the normal manner. What must be considered is whether the system WPARs that share a single installation will or will not interfere with each other's operation.

For software to function correctly in a shared-installation environment, the software package must be split into shareable and non-shareable files:

- ▶ Shareable files (such as executable code and message catalogs) must be installed into the shared global file systems that are read-only to all system WPARs.
- ▶ Non-shareable files (such as configuration and runtime-modifiable files) must be installed into the file systems that are writable to individual WPARs. This configuration allows multiple WPARs to share a single installation, yet still have unique configuration and runtime data.

In addition to splitting the software package, the software installation process must include a synchronization step to install non-shareable files into system WPARs. To accomplish this task, the application must provide a means to encapsulate the non-shareable files within the shared global file systems so that the non-shared files can be extracted into the WPAR by some means. For example, if a vendor creates a custom-installation system that delivers files into `/usr` and `/`, then the files that are delivered into `/` must be archived within `/usr` and then extracted into `/` using some vendor-provided mechanism. This action can occur automatically the first time that the application is started or configured.

Finally, the software update process must work so that the shareable and non-shareable files stay synchronized. If the shared files in the global AIX instance are updated to a certain fix level, then the non-shared files in individual WPARs also must be updated to the same level. Either the update process discovers all the system WPARs that must be updated or, at start time, the application detects the out-of-synchronization condition and applies the update. Some software products manage to never change their non-shareable files in their update process, so they do not need any special handling for updates.

This type of installation sometimes takes a little effort on the part of the application, but it allows you to get the most value from using WPARs. If there is a need to run the same version of the software in several WPARs, this type of installation provides the following benefits:

- ▶ It increases administrative efficiency by reducing the number of application instances that users must maintain. The administrator saves time in application-maintenance tasks, such as applying fixes and performing backups and migrations.
- ▶ It allows users to quickly deploy multiple instances of the same application, each in its own secure and isolated environment. It can take only a matter of minutes to create and start a WPAR to run a shared installation of the application.
- ▶ By sharing one AIX or application image among multiple WPARs, the memory resource usage is reduced because only one copy of the application image is in real memory.

For more information about WPAR, see *WPAR concepts*, available here:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.wpar/wpar-overview.htm>

Information about the topic of operating system-specific optimizations, from the IBM i and Linux perspectives, is available here:

- ▶ 5.3, “IBM i operating system-specific optimizations” on page 110 (*IBM i*)
- ▶ 6.3, “Linux operating system-specific optimizations” on page 126 (*Linux*)

4.4 AIX preferred practices

This section describes AIX preferred practices, and includes three subsections:

- ▶ AIX preferred practices that are applicable to all Power Systems generations
- ▶ AIX preferred practices that are applicable to POWER7 and POWER8 systems

4.4.1 AIX preferred practices that are applicable to all Power Systems generations

Preferred practices for the installation and configuration of all Power Systems generations are:

- ▶ If this server is a VIO Server, then run the VIO Performance Advisor on the VIO Server. Instructions are available for *Virtual I/O Server Advisor* at this website:

<http://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/Power%20Systems/page/VIOS%20Advisor>

For more information, see “VIO Performance Advisor” on page 201.

- ▶ For logical partitions (LPARs) with Java applications, run and evaluate the output from the Java Performance Advisor, which can be run on POWER5 and POWER6, to determine if there is an existing issue before you migrate to POWER7. Instructions are available for *Java Performance Advisor (JPA)* at this website:

[https://www.ibm.com/developerworks/community/wikis/home/wiki/Power%20Systems/page/Java%20Performance%20Advisor%20\(JPA\)](https://www.ibm.com/developerworks/community/wikis/home/wiki/Power%20Systems/page/Java%20Performance%20Advisor%20(JPA))

For more information, see “Java Performance Advisor” on page 203.

- ▶ For virtualized environments, you can also use the IBM PowerVM Virtualization Performance Advisor. Instructions for the IBM *PowerVM Virtualization Performance Advisor* are available here:
<https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/Power%20Systems/page/PowerVM%20Virtualization%20Performance%20Advisor>
For more information, see “Virtualization Performance Advisor” on page 202.
- ▶ The number of online virtual CPUs of a single LPAR cannot exceed the number of active CPUs in a pool. See the output of `lparstat -i` from the LPAR to see the values for online virtual CPUs and active CPUs in pool.
- ▶ IBM maintains a strong focus on the quality and reliability of Power Systems servers. To maintain this reliability, the currency of microcode levels on your systems is critical. Therefore, apply the latest Power Systems firmware and management console levels as soon as possible. These service pack updates contain a collective number of High Impact or PERvasive (HIPER) fixes that continue to provide you with the system availability you expect from IBM Power Systems.
- ▶ When you install firmware from the HMC, avoid the **do not auto accept** option. Selecting this advanced option can cause firmware installation problems.
- ▶ Subscribe to My Notifications to provide you with customizable communications that contain important news, new or updated support content, such as publications, hints, and tips, technical notes, product flashes (alerts), downloads, and drivers.

4.4.2 AIX preferred practices that are applicable to POWER7 and POWER8

The following are the AIX preferred practices that are applicable to POWER7 and POWER8.

Preferred practices for installation and configuration

Preferred practices for installation and configuration are:

- ▶ To ensure that your system conforms to the minimum requirements, see Chapter 3, “The POWER Hypervisor” on page 51 and the references that are provided for that chapter (see 4.5, “Related publications” on page 98).
- ▶ Review the *POWER7 Virtualization Best Practice Guide*, available here:
https://www.ibm.com/developerworks/wikis/download/attachments/53871915/P7_virtualization_bestpractice.doc?version=1
- ▶ For POWER7 and POWER7+ systems, review the discussion about the Active System Optimizer/Dynamic System Optimizer in SG248079 for identifying if those optimizations would be useful. Visit:
<http://www.redbooks.ibm.com/abstracts/sg248079.html>

For more information about this topic, see 4.5, “Related publications” on page 98.

4.5 Related publications

The publications that are listed in this section are considered suitable for a more detailed discussion of the topics that are covered in this chapter:

- ▶ *1 TB Segment Aliasing*, found here:

http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.prftungd/doc/prftungd/1TB_segment_aliasing.htm

- ▶ *AIX 64-bit Performance in Focus*, SG24-5103

- ▶ *AIX dscr_ctl API sample code*, found here:

<https://www.power.org/documentation/performance-guide-for-hpc-applications-on-ibm-power-755-system/> (registration required)

- ▶ *AIX Version 7.1 Release Notes*, found here:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.ntl/RELNOTES/GI11-9815-00.htm>

See the section, *The dscrctl command*.

- ▶ *Application configuration for large pages*, found here:

http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.prftungd/doc/prftungd/config_apps_large_pages.htm

- ▶ *AIX Linking and Loading Mechanisms*, found here:

http://download.boulder.ibm.com/ibmdl/pub/software/dw/aix/es-aix_ll.pdf

- ▶ *Efficient I/O event polling through the pollset interface on AIX*, found here:

<http://www.ibm.com/developerworks/aix/library/au-pollset/index.html>

- ▶ *Exclusive use processor resource sets*, found here:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.baseadm/doc/baseadmdita/excluseprocres.htm>

- ▶ *execrset command*, found here:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.cmds/doc/aixcmds2/execrset.htm>

- ▶ *General Programming Concepts: Writing and Debugging Programs*, found here:

http://publib16.boulder.ibm.com/doc_link/en_US/a_doc_lib/aixprggd/genprogc/understanding_mem_mapping.htm

- ▶ *IBM AIX Version 7.1 Differences Guide*, SG24-7910

See section 1.2, “Improved performance using 1 TB segments”

- ▶ *load and loadAndInit Subroutines*, found here:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.basetechref/doc/basetrf1/load.htm>

- ▶ *sync (Synchronize) or dcs (Data Cache Synchronize) instruction*, including information about **sync** and **lwsync** (lightweight sync), found here:

http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.assem/doc/aclangref/idalangref_sync_dcs_instrs.htm

- ▶ *mkrset Command*, found here:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.cmds/doc/aixcmds3/mkrset.htm>

- ▶ *Multiprocessing*, found here:
http://publib.boulder.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.prftungd/doc/prftungd/intro_multitproc.htm
- ▶ *Oracle Database and 1 TB Segment Aliasing*, found here:
<http://www.ibm.com/support/techdocs/atsmastr.nsf/WebIndex/TD105761>
- ▶ *pollset_create, pollset_ctl, pollset_destroy, pollset_poll, and pollset_query Subroutines*, found here:
<http://publib.boulder.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.basetechref/doc/basetrf1/pollset.htm>
- ▶ *The Performance of Runtime Data Cache Prefetching in a Dynamic Optimization System*, found here:
<http://www.microarch.org/micro36/html/pdf/lu-PerformanceRuntimeData.pdf>
- ▶ *POWER6 Decimal Floating Point (DFP)*, found here:
<http://www.ibm.com/developerworks/wikis/display/WikiPtype/Decimal+Floating+Point>
- ▶ *POWER7 Virtualization Best Practice Guide*, found here:
https://www.ibm.com/developerworks/wikis/download/attachments/53871915/P7_virtualization_bestpractice.doc?version=1
- ▶ *ra_attach Subroutine*, found here:
http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.basetechref%2Fdoc%2Fbasetrf2%2Fra_attach.htm
- ▶ *Shared library memory footprints on AIX 5L*, found here:
http://www.ibm.com/developerworks/aix/library/au-slib_memory/index.html
- ▶ *Simultaneous multithreading*, found here:
<http://publib.boulder.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.genprogc/doc/genprogc/smt.htm>
- ▶ *splat Command*, found here:
<http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.cmds/doc/aixcmds5/splat.htm>
- ▶ *trace Daemon*, found here:
<http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.cmds/doc/aixcmds5/trace.htm>
- ▶ *thread_post Subroutine*, found here:
http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.basetechref/doc/basetrf2/thread_post.htm
- ▶ *thread_post_many Subroutine*, found here:
http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.basetechref/doc/basetrf2/thread_post_many.htm
- ▶ *thread_wait Subroutine*, found here:
http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.basetechref/doc/basetrf2/thread_wait.htm

- *Thread environment variables*, found here:

https://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.prftungd/doc/prftungd/thread_env_vars.htm

- *Power ISA Version 2.07*, found here:

<https://www.power.org/documentation/power-isa-version-2-07/>

See the following sections:

- Section 3.1: Program Priority Registers
- Section 3.2: “or” Instruction
- Section 4.3.4: Program Priority Register
- Section 4.4.3: OR Instruction
- Section 5.3.4: Program Priority Register
- Section 5.4.2: OR Instruction
- Book I – 4 Floating Point Facility
- Book I – 5 Decimal Floating Point
- Book I – 6 Vector Facility
- Book I – 7 Vector-Scalar Floating Point Operations (VSX)
- Book I – Chapter 5 Decimal Floating-Point.
- Book II – 4.2 Data Stream Control Register
- Book II – 4.3.2 Data Cache Instructions
- Book II – 4.4 Synchronization Instructions
- Book II – A.2 Load and Reserve Mnemonics
- Book II – A.3 Synchronize Mnemonics
- Book II – Appendix B. Programming Examples for Sharing Storage
- Book III – 5.7 Storage Addressing



IBM i

This chapter describes the optimization and tuning of POWER8 and other Power processor-based servers running the IBM i operating system. It covers the following topics:

- ▶ 5.1, “Introduction” on page 102
- ▶ 5.2, “Using Power features with IBM i” on page 102
- ▶ 5.3, “IBM i operating system-specific optimizations” on page 104
- ▶ 5.4, “Related publications” on page 106

5.1 Introduction

IBM i provides an operating environment that emphasizes integration, security, and ease of use.

5.2 Using Power features with IBM i

The operating system and most applications for IBM i are built on a Technology Independent Machine Interface (TIMI) that isolates programs from differences in processor architectures, and allows the system to *automatically capitalize* on many new Power architecture features without changes to existing programs. For example, TIMI allows a program to use decimal floating point (DFP) on POWER5 (without special hardware support), and that same program automatically uses hardware support for DFP on POWER6, POWER7, and on POWER8 systems.

IBM Portable Application Solutions Environment for i (PASE for i) is a part of IBM i that allows some AIX application binaries to run on IBM i with little or no changes, so many optimizations described for AIX are applicable to PASE for i.

5.2.1 Multi-core and multi-thread

Operating system enablement usage of multi-core and multi-thread technology varies by operating system and release. Table 5-1 shows the maximum processor cores, threads, and SMT level for a (single) logical partition running IBM i. Customers who need more capacity can contact IBM to request more processor cores than are shown in Table 5-1. IBM Systems Lab Services works with clients to determine whether IBM i can support the customer workload in a partition with a larger number of cores.

Table 5-1 Maximum processor cores, threads, and SMT level for a (single) logical partition running IBM i

Release	POWER6	POWER7	POWER8
IBM i 6.1	32 Cores / 64 Threads / SMT2	Not supported	Not supported
IBM i 6.1.1	32 Cores / 64 Threads / SMT2	32 Cores / 128 Threads / SMT4	Not supported
IBM i 7.1 TR8	32 Cores / 64 Threads / SMT2	32 Cores / 128 Threads / SMT4	32 Cores / 256 Threads / SMT8
IBM i 7.2	32 Cores / 64 Threads / SMT2	32 Cores / 128 Threads / SMT4	48 Cores / 384 Threads / SMT8

Further information about this topic, from the processor and OS perspectives, is available here:

- ▶ 2.2.1, “Multi-core and multi-thread” on page 23 (*processor*)
- ▶ 4.2.1, “Multi-core and multi-thread” on page 64 (*AIX*)
- ▶ 6.2.1, “Multi-core and multi-thread” on page 108 (*Linux*)

SMT

SMT is a feature of the Power architecture and is described in “SMT” on page 25.

SMT dispatch control

IBM i 7.2 adds a job attribute named *Processor Resources Priority* (PRCRSCPTY) to influence how threads for the job are dispatched. The PRCRSCPTY attribute can request that the system isolate threads for the job on processors that are running fewer threads concurrently, or that the system run threads for the job on processors that are running as many concurrent threads as possible.

Information about the topic of SMT, from the processor and OS perspectives, is available here:

- ▶ “SMT” on page 22 (*processor*)
- ▶ “SMT” on page 61 (*AIX*)
- ▶ “SMT” on page 107 (*Linux*)

5.2.2 Multipage size support on IBM i

Most of IBM i uses 4 KB pages, but select system functions automatically use 64 KB pages. Applications running on IBM i 6.1 or later can create shared memory objects that use 64 KB pages (typically using `shmctl` with `SHM_PAGESIZE`). IBM technology for Java programs running on IBM i 6.1 or later can use 64 KB pages for Java heap. PASE for i programs running on IBM i 7.1 or later automatically use 64 KB pages for shared library text and data, and can request 64 KB pages for program text, stack, and data.

IBM Power Systems firmware does not support 64 KB pages for all configurations. For example, 64 KB pages are not available in a logical partition that is configured for Active Memory Sharing (AMS).

Information about this topic, from the processor and OS perspectives, is available here:

- ▶ 2.2.2, “Multipage size support: Page sizes (4 KB, 64 KB, 16 MB, and 16 GB)” on page 27 (*processor*)
- ▶ 4.2.2, “Multipage size support on AIX” on page 74
- ▶ 6.2.2, “Multipage size support on Linux” on page 113

5.2.3 Vector Scalar eXtension (VSX)

IBM i 7.2 automatically uses POWER8 vector instructions to improve the performance of some cryptographic operations. PASE for i applications running on IBM i 7.2 on POWER7 or newer processors can use VSX.

Information about the topic of VSX, from the processor, OS, and compiler perspectives, is available here:

- ▶ 2.2.5, “Vector Scalar eXtension (VSX)” on page 40 (*processor*)
- ▶ 4.2.5, “Vector Scalar eXtension (VSX)” on page 82 (*AIX*)
- ▶ 6.2.5, “Vector Scalar eXtension (VSX)” on page 120 (*Linux*)
- ▶ 7.3.2, “Compiler support for VSX” on page 145 (*XL and GCC compiler families*)

5.2.4 Decimal floating point

IBM i 6.1 and later provides support for DFP in select programming languages and in DB2 for i. DFP operations (outside of PASE for i) automatically use DFP instructions when running on POWER6 or newer processors, and use software support on older architectures. IBM i 7.2 improves the performance of many DFP operations (compared to prior releases) by the increased use of DFP instructions.

Information about this topic, from the processor and OS perspectives, is available here:

- ▶ 2.2.6, “Decimal floating point” on page 42 (*processor*)
- ▶ 4.2.6, “Decimal floating point (DFP)” on page 83 (*AIX*)
- ▶ 6.2.6, “Decimal floating point (DFP)” on page 120 (*Linux*)

5.3 IBM i operating system-specific optimizations

Here we describe optimization methods specific to IBM i.

5.3.1 IBM i advanced optimization techniques

Optimization methods specific to the creation of IBM i programs and service programs include the following capabilities¹:

- ▶ *8-byte pointers in C and C++ code*: The performance of C and C++ code that uses pointers can be improved when the code is compiled to use 8-byte pointers, rather than 16-byte pointers (default). To take full advantage of 8-byte pointers, specify STGM DL(*TERASPACE) and DTAMD L(*LLP64) when you compile code.
- ▶ *Program profiling*: Program profiling is an advanced optimization technique to reorder procedures, or code within procedures, and to direct code generation decisions in ILE programs and service programs based on statistical data gathered while running the program. The reordering can improve instruction cache utilization and reduce paging required by the program, thereby improving performance.
- ▶ *Argument optimization*: The Argument optimization parameter, with ARGOPT(*YES), is available with the CRTPGM and CRTSRVPGM commands to support advanced argument optimization, where an analysis across modules bound to the program is performed. In general, this improves the performance of most procedure calls within the program. Argument optimization is a technique for passing arguments (parameters) to ILE procedures to improve performance of call-intensive applications.
- ▶ *Interprocedural analysis*: Interprocedural analysis that is performed by the IPA(*YES) option on CRTPGM or CRTSRVPGM performs optimizations across function bodies in the entire program, during program creation. In particular, this occurs across the modules that are bound into the program and which were compiled with the MODCRTOPT(*KEEPILDTA) option. In contrast, intraprocedural is a mechanism for performing optimization for each function within a compilation unit, using only the information that is available for that function and compilation unit.
- ▶ *Licensed Internal Code Options (LICOPTs)*: LICOPTs are compiler options that are passed to the Licensed Internal Code to affect how code is generated or packaged. You can use some of the options to fine-tune the optimization of your code.

The TargetProcessorModel LICOPT instructs the translator to perform optimizations that are tuned for the specified processor model. Programs created with this option run on all supported hardware models, but generally run faster on the specified processor model. For IBM i 7.2, a TargetProcessorModel value can be specified so that the code should be tuned to optimally run on POWER8.

The CodeGenTarget LICOPT specifies the creation target model for a program or module object. The creation target model indicates the hardware features that the code generated for that object can use. For IBM i 7.2, a CodeGenTarget model of POWER8 can be specified.

CodeGenTarget features and their associated Power Systems hardware include these:

- CodeGenTarget features associated with POWER6 hardware:
 - A hardware decimal floating point unit
 - Efficient hardware support for ILE pointer handling

¹ ILE Concepts (SC41-5606),

<http://pic.dhe.ibm.com/infocenter/iseries/v7r1m0/topic/rzahg/rzahgileconcept.htm>

- CodeGenTarget features associated with POWER7 hardware:
 - A number of new instructions that may speed up certain computations, such as conversions between integer and floating-point values.
- CodeGenTarget features associated with POWER8 hardware:
 - New move instructions between floating point and general purpose registers

The TargetProcessorModel and CodeGenTarget LICOPTs are two of several factors, including Adaptive Code Generation, that determine the processor model to which code should be tuned and targeted when creating a module, changing a module or program, or re-creating a module or program. The default behavior is to use all features available on the current machine. See *Related publications* at the end of this chapter for a link to *ILE Concepts (SC41-5606)*.

- *Adaptive Code Generation (ACG)*: ACG allows you to take advantage of all of the processor features on your systems, regardless of whether those features are present on other system models that are supported by the same release. Furthermore, programs can be moved from one system model to another and continue to run correctly, even if the new machine does not have all of the processor features that were available on the original machine. The technology for achieving this is ACG. ACG can work without user intervention in most scenarios. However, if you build and distribute software to run on a variety of system models, you might want to exercise some control over which processor features are used by ACG.

The first time a program object is activated on a system to which it is moved, the system performs a compatibility check to ensure that your program does not use any features that are unavailable on your system. If the program requires any processor feature that is not supported by the system to which it was moved, then the system automatically calls the optimizing translator to convert the program to be compatible. Options that are associated with restoring objects exist to cause incompatible module and program objects that are restored to be immediately converted, rather than on the first activation.

Additional information about these optimizations in an IBM i environment can be found here:

- *ILE Concepts*, available here:
<http://pic.dhe.ibm.com/infocenter/iserics/v7r1m0/topic/rzahg/rzahgileconcept.htm>

In particular, see the chapter about advanced optimization techniques.

5.3.2 Performance management on IBM i

For links to general performance resources, performance education resources, performance papers, and articles for IBM i, see the following reference:

- *Performance management on IBM i*, available here:
<http://www.ibm.com/systems/power/software/i/management/performance/resources.html>

For a basic understanding of IBM i on Power Systems performance concepts, workloads and benchmarks on Power Systems, capacity planning, performance monitoring and analysis, frequently asked questions, and guidelines addressing common performance issues, see the following reference:

- *IBM i on Power - Performance FAQ*, available here:
http://www-01.ibm.com/common/ssi/cgi-bin/ssialias?subtype=WH&infotype=SA&appname=STGE_PO_PO_USEN&htmlfid=POW03102USEN&attachment=POW03102USEN.PDF

Information about the topic of operating system-specific optimizations, from the AIX and Linux perspectives, is available here:

- ▶ 4.3, “AIX operating system-specific optimizations” on page 86 (*AIX*)
- ▶ 6.3, “Linux operating system-specific optimizations” on page 123 (*Linux*)

5.4 Related publications

- ▶ *ILE Concepts (SC41-5606)*, available here:
<http://pic.dhe.ibm.com/infocenter/iserics/v7r1m0/topic/rzahg/rzahgileconcept.htm>
- ▶ *Advanced Optimization Techniques*, available here:
<http://pic.dhe.ibm.com/infocenter/iserics/v7r1m0/topic/ilec/sc415606206.htm>
- ▶ *Performance management on IBM i*, available here:
<http://www.ibm.com/systems/power/software/i/management/performance/resources.html>
- ▶ *IBM i on Power - Performance FAQ*, available here:
http://www-01.ibm.com/common/ssi/cgi-bin/ssialias?subtype=WH&infotype=SA&appname=STGE_PO_PO_USEN&htmlfid=POW03102USEN&attachment=POW03102USEN.PDF



Linux

This chapter describes the optimization and tuning of the POWER8 and other Power processor-based servers running the Linux operating system. It covers the following topics:

- ▶ 6.1, “Introduction” on page 108
- ▶ 6.2, “Using Power features with Linux” on page 108
- ▶ 6.3, “Linux operating system-specific optimizations” on page 123
- ▶ 6.4, “Related publications” on page 132

6.1 Introduction

When you work with POWER7, POWER7+, or POWER8 processor-based servers, systems, and solutions, a solid choice for running enterprise-level workloads is Linux. Red Hat Enterprise Linux (RHEL) and SUSE Linux Enterprise Server (SLES) provide operating systems that are optimized and targeted for the Power Architecture. These operating systems run natively on the Power Architecture and are designed to take full advantage of the specialized features of Power Systems.

Both RHEL and SLES provide the tools, kernel support, optimized compilers, and tuned libraries for POWER7 and POWER8 systems. The Linux distributions provide for excellent performance, and more application and customer-specific tuning approaches are available. IBM provides a number of value-add packages, tools, and extensions that provide for more tunings, optimizations, and products for the best possible performance on POWER8. The typical Linux open source performance tools that Linux users are comfortable with are available on the IBM PowerLinux systems.

The IBM PowerLinux Tools repository enables the use of standard Linux package management tools (such as yum and zypper) to provide easy access to IBM preferred tools:

- ▶ IBM PowerLinux hardware diagnostic aids and productivity tools
- ▶ IBM Software Development Toolkit for PowerLinux servers
- ▶ IBM Advance Toolchain for PowerLinux servers

The IBM PowerLinux Tools repository is available here:

<http://www.ibm.com/support/customer/sas/f/1opdiags/yum.html>

The Linux distributions are enabled to run on the broad range of IBM Power offerings, from low-cost PowerLinux servers and Flex System nodes, up through the largest IBM Power 770 and Power 795 servers. Linux on Power supports small virtualized Micro-Partitioning partitions up through large dedicated partitions containing all of the resources of a high-end server. For KVM-enabled PowerLinux systems, the Linux distributions are enabled to run in the KVM guests on the PowerLinux servers.

IBM premier products, such as IBM XL compilers, IBM Java products, IBM WebSphere, and IBM DB2 database products, all provide Power optimized support with the RHEL and SLES operating systems.

For more information about this topic, see 6.4, “Related publications” on page 132.

6.2 Using Power features with Linux

Some of the significant features of POWER with POWER7 and POWER8 extensions in a Linux environment are described in this section.

6.2.1 Multi-core and multi-thread

Operating system enablement usage of multi-core and multi-thread technology varies by operating system and release. Table 6-1 shows the maximum processor cores and threads for a (single) logical partition running Linux.

Table 6-1 Multi-thread per core features by single LPAR scaling

Single LPAR scaling	Linux release
32-core/32-thread	RHEL 5/6, SLES 10/11
64-core/128-thread	RHEL 5/6, SLES 10/11
64-core/256-thread	RHEL 6, SLES 11sp1
256-core/1024-thread	RHEL 6, SLES 11sp1

Information about multi-thread per core features by POWER generation is available in Table 2-1 on page 24.

Further information about this topic, from the processor and OS perspectives, is available here:

- ▶ 2.2.1, “Multi-core and multi-thread” on page 23 (*processor*)
- ▶ 4.2.1, “Multi-core and multi-thread” on page 64 (*AIX*)
- ▶ 5.2.1, “Multi-core and multi-thread” on page 102 (*IBM i*)

Simultaneous multithreading (SMT)

Simultaneous multithreading (SMT) is a feature of the Power architecture and is described in “SMT” on page 25.

On a POWER8 system, with a properly enabled Linux distribution, or distro, the Linux operating system supports up to eight hardware threads per core (SMT=8).

With the POWER8 processor cores, the SMT hardware threads are designed to be more equal in the execution implementation, which allows the system to support flexible SMT scheduling and management.

Application throughput and SMT scaling from SMT=1 to SMT=2, to SMT=4, and to SMT=8 is highly application dependent. With additional hardware threads available for scheduling, the ability of the processor cores to switch from a waiting (stalled) hardware thread to another thread that is ready for processing can improve overall system effectiveness and throughput.

High SMT modes are best for maximizing total system throughput, while lower SMT modes might be appropriate for high performance threads and low latency applications. For code with low levels of instruction-level parallelism (often seen in Java code, for example), high SMT modes are generally preferred.

Information about the topic of SMT, from the processor and OS perspectives, is available here:

- ▶ “SMT” on page 25 (*processor*)
- ▶ “Simultaneous Multithreading (SMT)” on page 65 (*AIX*)
- ▶ “SMT” on page 102 (*IBM i*)

Boot-time enablement of SMT

When booting a Linux distro that supports SMT=8, SMT=8 is the default boot mode. The system can be booted to SMT off or the default SMT on mode by adding an **smt-enabled=off** or **smt-enabled=on** kernel parameter to the append line of the bootloader file.

Dynamically selecting different SMT modes

Linux enables Power SMT capabilities. By default, the system runs at the highest SMT level.

Changing SMT settings remains a dynamic (runtime) option in the operating system. The **ppc64_cpu** command is provided in the `powerpc_utils` package. Running this command requires root access. The **ppc64_cpu** command can be used to force the system kernel to use lower SMT levels (ST, SMT2, or SMT4 mode). Here is an example:

- ▶ **ppc64_cpu --smt=1** sets the SMT mode to ST
- ▶ **ppc64_cpu --smt** shows the current SMT mode

POWER8 systems support up to 8 SMT hardware threads per core. The **ppc64_cpu** command can specify hardware threads from a single thread per core, 2 threads, 4 threads, or 8 threads.

When using the **ppc64_cpu** command to control SMT settings, the normal Linux approach of *holes in the CPU numbering* continues as it was in previous POWER generations, such as POWER7.

In different POWER8 SMT modes, CPUs are numbered as follows:

SMT=8:	0,1,2,3,4,5,6,7,	8,9,10,11,12,13,14,15,	16,17,18,19,20,21,22,23,	and so on
SMT=4:	0,1,2,3,	8,9,10,11,	16,17,18,19,	and so on
SMT=2:	0,1,	8,9,	16,17,	and so on
SMT=1:	0,	8,	16,	and so on

The setaffinity application programming interface (API) allows processes and threads to have affinity to specific logical processors. See “Affinitization and binding” on page 111. Because POWER8 supports running up to 8 threads per core, the CPU numbering is different than in POWER7, which only supported up to 4 threads per core. Therefore, an application that specifically binds processes to threads will need to be aware of the new CPU numbering to ensure the binding is correct, because there are now more threads available for each core.

For more information about this topic, see 6.4, “Related publications” on page 132.

Querying the SMT setting

The command for querying the SMT setting is `ppc64_cpu --smt`. A programmable API is not available.

SMT priorities

SMT priorities in the Power hardware are introduced in “SMT priorities” on page 25. Linux supports selecting SMT priorities using the Priority Nop mechanism or by writing to the PPR, as described in that section.

The current GLIBC (from version 2.16) and forward provide the system header `sys/platform/ppc.h` which contains a wrapper for setting the PPR using the Priority Nop mechanism. See Example 6-1.

Example 6-1 6-X GLIBC PPR set functions

```
void __ppc_set_ppr_med (void)
void __ppc_set_ppr_med_low (void)
void __ppc_set_ppr_low (void)
```

RHEL6 and SLES11p1 do not provide this header; however, it is supported on AT 6.0 and later.

Where to use

SMT thread priority can be used to improve the performance of a workload by lowering the SMT thread priority that is being used on an SMT thread that is running a particular process-thread in these circumstances:

- ▶ The thread is waiting on a lock.
- ▶ The thread is waiting on an event, such as the completion of an I/O event.

Alternatively, process-threads that are performance sensitive can maximize their performance by ensuring that the SMT thread priority level is set to an elevated level.

Information about the topic of SMT priorities, from the processor and OS perspectives, is available here:

- ▶ “SMT priorities” on page 25 (*processor*)
- ▶ “SMT priorities” on page 66 (*AIX*)

Affinitization and binding

Affinity performance effects are explained in “The POWER8 processor and affinity performance effects” on page 14. Establishing good affinity is accomplished by understanding the placement of a partition on the underlying cores and memory of a Power system, and then by using operating system facilities to bind application threads to run on specific hardware threads or cores.

The `numactl --hardware` command shows the relative positioning of the underlying cores and memory, if that information is available from the hypervisor or firmware. In the case of PowerVM shared pools, this information cannot be directly mapped to the underlying cores and memory.

Flexible SMT support

On POWER7 and POWER7+, there is a correlation between the hardware thread number (0-3) and the hardware resources within the processor. Matching the thread numbers to the number of active threads has been advised for optimum performance. For example, if only one thread is active, it should be hardware thread 0. If two threads are active, they should be hardware threads 0 and 1. The Linux operating system automatically shifts the threads to those modes.

On POWER8, any process or thread can run in any SMT mode. The processor balances the processor core resources according to the number of active hardware threads. There is no need to match the application thread numbers with the number of active hardware threads. Hardware threads on the POWER8 processor have equal weight, unlike the hardware threads under POWER7. Therefore, as an example, a single process running on thread 7 would run just as fast as running on thread 0, presuming nothing else is on the other hardware threads for that processor core.

Linux scheduler

The Linux Completely Fair Scheduler (CFS) handles load balancing across CPUs and uses scheduler modules to make policy decisions. CFS works with multi-core, multi-thread processors and will balance tasks across real processors. CFS also groups and tunes related tasks together.

The Linux topology considers physical packages, threads, siblings, and cores. The CFS scheduler domains help to determine load balancing. The base domain contains all sibling threads of the physical CPU; the next parent domain contains all physical CPUs; and the next parent domain takes NUMA nodes into consideration.

Due to the specific asymmetrical thread ordering of POWER7, special Linux scheduler modifications had to be added for the POWER7 CPU type. With POWER8, this logic is no longer needed, because any of the SMT8 threads can act as the primary thread by design.

This means the number of threads that are active in the core at one time determines the dynamic SMT mode (for example, from a performance perspective, thread 0 can be the same as thread 7). Idle threads should be napping or in a deeper sleep if they are idle for a period of time.

CPUsets, cgroups, scheduler domains

It is possible to target (and limit) processes for a specific set of CPUs or cores. This can provide Linux applications with more fine-grained control of the cores and characteristics of application process and thread requirements.

taskset

Use the **taskset** command to retrieve, set, and verify the CPU affinity information of a process that running.

numactl

Similar to the **taskset** command, use the **numactl** command to retrieve, set, and verify the CPU affinity information of a process that running. The **numactl** command, however, provides additional performance information about local memory allocation.

Using setaffinity to bind to specific logical processors

The setaffinity API allows processes and threads to have affinity to specific logical processors. The number and numbering of logical processors is a product of the number of processor cores (in the partition) and the SMT capability of the machine (eight-way SMT for POWER8).

Information about the topic of affinitization and binding, from the processor and OS perspectives, is available here:

- ▶ “Affinitization and binding to hardware threads” on page 26 (*processor*)
- ▶ “Affinitization and binding” on page 66 (*ALX*)

Hybrid thread and core

Linux provides facilities to customize SMT characteristics of CPUs running within a partition.

SMT can be enabled or disabled at boot time as described in “Boot-time enablement of SMT” on page 109. SMT modes can be dynamically controlled for a whole partition, as described in “Dynamically selecting different SMT modes” on page 109.

In Linux, each CPU is associated with a processor core hardware thread.

When there is no work to be done on a CPU, the scheduler goes into the idle loop, and Linux calls into the hypervisor to report that the CPU is truly idle. The kernel-to-hypervisor interface is defined in the Power Architecture Platform Reference (PAPR) at <http://power.org>. In this case, it is the H_CEDCE hypervisor call.

Information about this topic, from the processor and OS perspectives, is available here:

- ▶ “Hybrid thread and core” on page 27 (*processor*)
- ▶ “Hybrid thread and core” on page 71 (*ALX*)

6.2.2 Multipage size support on Linux

On Power Systems running Linux, the default page size is 64 KB, so most, but not all, applications are expected to see a performance benefit from this default. There are cases in which an application uses many small files, which can mean that each file is loaded into a 64 KB page, resulting in poor memory utilization.

Support for 16 MB pages (hugepages in Linux terminology) is available through various mechanisms and is typically used for databases, Java engines, and high-performance computing (HPC) applications. The `libhugetlbfs` package is available in Linux distributions, and using this package gives you the most benefit from 16 MB pages.

Information about this topic, from the processor and OS perspectives, is available here:

- ▶ 2.2.2, “Multipage size support: Page sizes (4 KB, 64 KB, 16 MB, and 16 GB)” on page 27 (*processor*)
- ▶ 4.2.2, “Multipage size support on AIX” on page 74
- ▶ 5.2.2, “Multipage size support on IBM i” on page 103

6.2.3 Efficient use of cache

Operating system facilities for controlling hardware prefetching are described in this section.

Controlling DSCR under Linux

DSCR settings on Linux are controlled with the `ppc64_cpu` command. Controlling DSCR settings for an application is generally considered advanced and specific tuning.

Currently, setting the DSCR value is a cross-LPAR setting.

Information about the efficient use of cache, from the processor and OS perspectives, is available here:

- ▶ 2.2.3, “Efficient use of cache and memory” on page 28 (*processor*)
- ▶ 4.2.3, “Efficient use of cache” on page 78 (*AIX*)

For more information about this topic, see 6.4, “Related publications” on page 132.

6.2.4 Transactional memory (TM)

Transactional memory (TM) is a new approach that simplifies concurrent programming, specifically in the accessing of shared data across multiple threads. Previous to TM, accesses to shared data were synchronized by the use of locks. Threaded code that needed access to shared data first has to acquire the data lock, then access the shared data, and then release the lock. On many systems, acquiring locks can be expensive, making accessing shared data vastly more expensive than accessing non-shared data. This additional locking can be especially burdensome when the shared data has low contention between the multiple threads.

Using TM, shared data accesses are placed into blocks of code called *transactions*. When using hardware TM (HTM), these transactions are executed without locking, and the results are seen by other threads atomically.

The POWER8 processor supports the HTM instructions defined in the *Power ISA Version 2.07* documentation, located here:

<https://www.power.org/documentation/power-isa-version-2-07/>

Users have three options when writing code to make use of the POWER8 HTM features:

1. The first option of HTM use is through the low level GCC built-in functions, which are enabled with the GCC `-mcpu=power8` or `-mhtm` compiler options. The HTM built-in functions return `true` or `false`, depending on their success. The arguments to the HTM built-in functions match exactly the type and order of the associated hardware instruction operands, as shown in Example 6-2.

Example 6-2 GCC HTM built-in functions

```
unsigned int __builtin_tbegin (unsigned int)
unsigned int __builtin_tend (unsigned int)

unsigned int __builtin_tabort (unsigned int)
unsigned int __builtin_tabortdc (unsigned int, unsigned int, unsigned int)
unsigned int __builtin_tabortdci (unsigned int, unsigned int, int)
unsigned int __builtin_tabortwc (unsigned int, unsigned int, unsigned int)
unsigned int __builtin_tabortwci (unsigned int, unsigned int, int)

unsigned int __builtin_tcheck (unsigned int)
unsigned int __builtin_treclaim (unsigned int)
unsigned int __builtin_trechktpt (void)
unsigned int __builtin_tsr (unsigned int)

unsigned long __builtin_get_texasr (void)
unsigned long __builtin_get_texasru (void)
unsigned long __builtin_get_tfhar (void)
unsigned long __builtin_get_tfiar (void)

void __builtin_set_texasr (unsigned long);
void __builtin_set_texasru (unsigned long);
void __builtin_set_tfhar (unsigned long);
void __builtin_set_tfiar (unsigned long);
```

In addition to Example 6-2, HTM built-in functions, we have added built-in functions for some common extended mnemonics of the HTM instructions, as shown in Example 6-3.

Example 6-3 GCC HTM built-in functions for extended mnemonics

```
unsigned int __builtin_tendall (void)
unsigned int __builtin_tresume (void)
unsigned int __builtin_tsuspend (void)
```

Common usage of these HTM built-in functions might produce results similar to those shown in Example 6-4.

Example 6-4 Simple use of HTM built-in functions

```
#include <htmintrin.h>
if (__builtin_tbegin (0))
{
    /* Transaction State Initiated. */
    if (is_locked (lock))
        __builtin_tabort (0);
    a = b + c;
    __builtin_tend (0);
}
else
{
    /* Transaction State Failed, Use Locks. */
    acquire_lock (lock);
    a = b + c;
    release_lock (lock);
}
```

A slightly more complicated example is shown in Example 6-5. Here, we attempt to retry the transaction a specific number of times before falling back to using locks.

Example 6-5 Complex use of HTM built-in functions

```
#include <htmintrin.h>
int num_retries = 10;
while (1)
{
    if (__builtin_tbegin (0))
    {
        /* Transaction State Initiated. */
        if (is_locked (lock))
            __builtin_tabort (0);
        a = b + c;
        __builtin_tend (0);
        break;
    }
    else
    {
        /* Transaction State Failed. Use locks if the transaction
           failure is "persistent" or we've tried too many times. */
        if (num_retries-- <= 0
            || _TEXASRU_FAILURE_PERSISTENT (__builtin_get_texasru ()))
        {
            acquire_lock (lock);
            a = b + c;
            release_lock (lock);
            break;
        }
    }
}
```

In some cases, it can be useful to know whether the code that is being executed is in transactional state or not. Unfortunately, that cannot be determined by analyzing the HTM Special Purpose Registers (SPRs). That specific information is only contained within the Machine State Register (MSR) Transaction State (TS) bits which are not accessible by user code. To allow access to that information, we have added one final built-in function and some associated macros to help the user to determine what the transaction state is at a particular point in their code.

```
unsigned int __builtin_ttest (void)
```

Usage of the built-in function and its associated macro might look like the code shown in Example 6-6.

Example 6-6 Determining Transaction State

```
#include <htmintrin.h>

unsigned char tx_state = __builtin_ttest ();

if (_HTM_STATE (tx_state) == _HTM_TRANSACTIONAL)
{
    /* Code to use in transactional state. */
}
else if (_HTM_STATE (tx_state) == _HTM_NONTRANSACTIONAL)
{
    /* Code to use in non-transactional state. */
}
else if (_HTM_STATE (tx_state) == _HTM_SUSPENDED)
{
    /* Code to use in transaction suspended state. */
}
```

2. A second option for using HTM is by using the slightly higher level inline functions that are common to both the GCC and the IBM XL compilers. These HTM built-ins are defined in the **htmxlintrin.h** header file and are also mostly common between POWER and IBM System z (with a few exceptions) and can be used to write code that can be compiled on POWER or System z using either the IBM XL or GCC compilers. See Example 6-7.

```
long __TM_simple_begin (void)
long __TM_begin (void* const TM_buff)
long __TM_end (void)
void __TM_abort (void)
void __TM_named_abort (unsigned char const code)
void __TM_resume (void)
void __TM_suspend (void)

long __TM_is_user_abort (void* const TM_buff)
long __TM_is_named_user_abort (void* const TM_buff, unsigned char *code)
long __TM_is_illegal (void* const TM_buff)
long __TM_is_footprint_exceeded (void* const TM_buff)
long __TM_nesting_depth (void* const TM_buff)
long __TM_is_nested_too_deep(void* const TM_buff)
long __TM_is_conflict(void* const TM_buff)
long __TM_is_failure_persistent(void* const TM_buff)
long __TM_failure_address(void* const TM_buff)
long long __TM_failure_code(void* const TM_buff)
```

Using these built-in functions, we can create a more portable version of the code in Example 6-5 on page 116, so that it will work on POWER and on System z, using either GCC or the XL compilers. This more portable version is shown in Example 6-8.

Example 6-8 Complex HTM usage using portable HTM intrinsics

```

#ifdef __GNUC__
# include <htmxlintrin.h>
#endif

int num_retries = 10;
TM_buff_type TM_buff;

while (1)
{
    unsigned char tx_status = __TM_begin (TM_buff);
    if (tx_status == _HTM_TBEGIN_STARTED)
    {
        /* Transaction State Initiated. */
        if (is_locked (lock))
            __TM_abort ();
        a = b + c;
        __TM_end ();
        break;
    }
    else
    {
        /* Transaction State Failed. Use locks if the transaction
           failure is "persistent" or we've tried too many times. */
        if (num_retries-- <= 0
#ifdef __powerpc__
            || __TM_is_failure_persistent (TM_buff)
#elif defined (__s390__)
            || __TM_is_failure_persistent (tx_status)
#endif
        )
        {
            acquire_lock (lock);
            a = b + c;
            release_lock (lock);
            break;
        }
    }
}

```

3. The third and most portable option uses a high level language interface which is implemented by GCC and the GNUTM Library (LIBITM).

<http://gcc.gnu.org/wiki/TransactionalMemory>

This high level language option is enabled using the **-fgnu-tm** option (**-mcpu=power8** and **-mhtm** are not needed), and it provides a common transactional model across multiple architectures and multiple compilers using the **__transaction_atomic {...}** language

construct. The LIBITM library which is included with the GCC compiler has the ability to determine, at runtime, whether it is executing on a processor that supports HTM instructions, and, if so, it utilizes them in executing the transaction. Otherwise, it automatically falls back to using software TM which relies on locks. LIBITM also has the ability to retry a transaction using HTM if the initial transaction **begin** failed, similar to the complicated example (Example 6-5 on page 116). An example of the third option that is equivalent to the complicated examples (Example 6-5 on page 116 and Example 6-8 on page 119) is simple and is shown in Example 6-9.

Example 6-9 GNU TM Library (LIBITM) Usage

```
__transaction_atomic
{
    a = b + c;
}
```

Support for the HTM built-in functions, the XL HTM built-in functions, and LIBITM support will be in an upcoming Free Software Foundation (FSF) version of GCC. However, it is also available in the GCC 4.8-based compiler that is shipped in Advance Toolchain (AT) version 7.0.

Information about the topic of TM, from the processor, OS, and compiler perspectives, is available here:

- ▶ 2.2.4, “Transactional memory (TM)” on page 37 (*processor*)
- ▶ 4.2.4, “Transactional memory (TM)” on page 81 (*AIX*)
- ▶ 7.3.5, “Transactional memory (TM)” on page 149 (*XL and GCC compiler families*)

6.2.5 Vector Scalar eXtension (VSX)

GCC makes an interface available for PowerPC processors to access built-in functions. See the documentation for the revision of the GCC compiler that you are using at:

<http://gcc.gnu.org/onlinedocs>

Information about the topic of VSX, from the processor, AIX, IBM i, and compiler perspectives, is available here:

- ▶ 2.2.5, “Vector Scalar eXtension (VSX)” on page 40 (*processor*)
- ▶ 4.2.5, “Vector Scalar eXtension (VSX)” on page 82 (*AIX*)
- ▶ 5.2.3, “Vector Scalar eXtension (VSX)” on page 103 (*IBM i*)
- ▶ 7.3.2, “Compiler support for VSX” on page 145 (*XL and GCC compiler families*)

6.2.6 Decimal floating point (DFP)

Decimal (base 10) data is widely used in commercial and financial applications. However, most computer systems have only binary (base two) arithmetic. There are two binary number systems in computers: integer (fixed-point) and floating point. Unfortunately, decimal calculations cannot be directly implemented with binary floating point. For example, the value 0.1 needs an infinitely recurring binary fraction, whereas a decimal number system can represent it exactly, as one tenth. So, using binary floating point cannot ensure that results are the same as those results using decimal arithmetic.

In general, decimal floating point (DFP) operations are emulated with binary fixed-point integers. Decimal numbers are traditionally held in a binary-coded decimal (BCD) format. Although BCD provides sufficient accuracy for decimal calculation, it imposes a heavy cost in performance, because it is usually implemented in software.

IBM POWER6, POWER7, and POWER8 processor-based systems provide hardware support for DFP arithmetic. These microprocessor cores include a DFP unit that provides acceleration for the DFP arithmetic. The IBM Power instruction set is expanded to include the following:

- Fifty-four new instructions were added to support the DFP unit architecture. DFP can provide a performance boost for applications that are using BCD calculations.

How to take advantage of DFP unit on POWER

You can take advantage of the DFP unit on POWER with the following features:¹

- Native DFP language support with a compiler

The C draft standard includes the following new data types (these are native data types, as are int, long, float, double, and so on):

<code>_Decimal32</code>	7 decimal digits of accuracy
<code>_Decimal64</code>	16 decimal digits of accuracy
<code>_Decimal128</code>	34 decimal digits of accuracy

Note: The `printf()` function uses new options to print these new data types:

- `_Decimal32` uses `%Hf`
- `_Decimal64` uses `%Df`
- `_Decimal128` uses `%DDf`

- The IBM XL C/C++ Compiler, release 9 or later for AIX and Linux, includes native DFP language support. Here is a list of compiler options for IBM XL compilers that are related to DFP:
 - `-qdfp`: Enables DFP support. This option makes the compiler recognize DFP literal suffixes, and the `_Decimal32`, `_Decimal64`, and `_Decimal128` keywords.
 - `-qfloat=dfpemulate`: Instructs the compiler to use calls to library functions to handle DFP computation, regardless of the architecture level. You might experience performance degradation when you use software emulation.
 - `-qfloat=nodfpemulate` (the default when the `-qarch` flag specifies POWER6, POWER7 or POWER8): Instructs the compiler to use DFP hardware instructions.
 - `-D__STDC_WANT_DEC_FP__`: Enables the referencing of DFP-defined symbols.
 - `-ldfp`: Enables the DFP functionality that is provided by the Advance Toolchain on Linux.

For hardware supported DFP, with `-qarch=pwr6`, `-qarch=pwr7`, or `-qarch=pwr8`, use the following command:

```
cc -qdfp
```

For software emulation of DFP (on earlier processor chips), use the following command:

```
cc -qdfp -qfloat=dfpemulate
```

¹ How to compile DFPAL?, available here: <http://speleotrove.com/decimal/dfpal/compile.html>

- The GCC compilers for Power Systems also include native DFP language support.

As of SLES 11 SP1 and RHEL 6, and in accord with the Institute of Electrical and Electronics Engineers (IEEE) 754R, DFP is fully integrated with compiler and run time (printf and DFP math) support. For older Linux distribution releases (RHEL 5/SLES 10 and earlier), you can use the freely available Advance Toolchain compiler and run time. The Advance Toolchain runtime libraries can also be integrated with recent XL (V9+) compilers for DFP exploitation.

The latest Advance Toolchain compiler and run times can be downloaded from the following website:

<ftp://ftp.unicamp.br/pub/linuxpatch/toolchain/at/>

Advance Toolchain is a self-contained toolchain that does not rely on the base system toolchain for operability. In fact, it is designed to coexist with the toolchain shipped with the operating system. You do not have to uninstall the regular GCC compilers that come with your Linux distribution to use the Advance Toolchain.

The latest Enterprise distributions and Advance Toolchain run time use the Linux CPU tune library capability to automatically select hardware DFP or software implementation library variants, which are based on the hardware platform.

Here is a list of GCC compiler options for Advance Toolchain that are related to DFP:

- **-mhard-dfp** (the default when **-mcpu=power6**, **-mcpu=power7** or **-mcpu=power8** is specified): Instructs the compiler to directly take advantage of DFP hardware instructions for decimal arithmetic.
 - **-mno-hard-dfp**: Instructs the compiler to use calls to library functions to handle DFP computation, regardless of the architecture level. If your application is dynamically linked to the libdfp variant and running on POWER6, POWER7, or POWER8 processors, then the run time automatically binds to the **libdfp** variant implemented with hardware DFP instructions. Otherwise, the software DFP library is used. You might experience performance degradation when you use software emulation.
 - **-D__STDC_WANT_DEC_FP__**: Enables the reference of DFP defined symbols.
 - **-ldfp**: Enables the DFP functionality that is provided by recent Linux Enterprise Distributions or the Advance Toolchain run time.
- Decimal Floating Point Library (**libdfp**) is an implementation of the joint efforts of the International Organization for Standardization and the International Electrotechnical Commission (ISO/IEC). ISO/IEC technical report ISO/IEC TR 24732² describes the C-Language library routines that are necessary to provide the C library runtime support for decimal floating point data types, as introduced in IEEE 754-2008, namely **_Decimal32**, **_Decimal64**, and **_Decimal128**.

The library provides functions, such as **sin** and **cos**, for the decimal types that are supported by GCC. Current development and documentation can be found at <https://github.com/libdfp/libdfp>, and RHEL6 and SLES11 provide this library as a supplementary extension. Advance Toolchain also ships with the library.

Determining if your applications are using DFP

The Linux **perf** tool is used for application monitoring. The **PM_MRK_DFU_FIN** performance counter event indicates that the Decimal Floating Point Unit finished a marked instruction.

² Information technology -- Programming languages, their environments and system software interfaces -- Extension for the programming language C to support decimal floating-point arithmetic
http://www.iso.org/iso/catalogue_detail.htm?csnumber=38842

To profile an application for PM_MRK_DFU_FIN samples, use **perf** to set the event name and sample count and run the application:

```
perf -e PM_MRK_DFU_FIN:1000 application
```

To view the results and see what symbols the event samples are associated with, use:

```
opreport --symbols
```

If you see this message, there were no samples found for the event specified when running the application:

```
opreport error: No sample file found
```

Information about this topic, from the processor and OS perspectives, is available here:

- ▶ 2.2.6, “Decimal floating point” on page 42 (*processor*)
- ▶ 4.2.6, “Decimal floating point (DFP)” on page 83 (*ALX*)
- ▶ 5.2.4, “Decimal floating point” on page 103 (*IBM i*)

For more information, see 6.4, “Related publications” on page 132.

6.2.7 Event-based branches

The event-based branching (EBB) facility is a new Power Architecture ISA 2.07 hardware facility, under [Category:Server], that generates event-based exceptions when a certain event criteria is met. Currently, ISA 2.07 (on POWER8 hardware) only defines one type of EBB: the performance monitoring unit (PMU) EBB. Following an EBB exception, the branch event status and control register (BESCR) tells which kind of event triggered the exception.

The EBB facility is a per-hardware-thread problem-state facility with access to the PMU and initialization under privileged-state. A problem-state application with direct access to the facility will register a callback function as an EBB handler by setting the handler address into the EBBHR register.

When a specified or requested PMU overflows, an exception is generated and, as a result, the problem-state application EBB handler is invoked. Execution continues in event-based exception context until the handler returns control to the address in the event-based branch return register (EBBRR) using the **rfebb** instruction.

There are interoperability considerations with the Power Architecture executable and linkable format (ELF) Application Binary Interface (ABI), and these can complicate the usage of this facility in order to remain ABI compliant. As a result, it is suggested that user applications utilize an API provided by **libpaf-ebb** that handles the ABI implications consistently and correctly and provides a handler by proxy.

Information about this topic, from the processor perspective, is available here:

- ▶ 2.2.12, “Event-based branches (or user-level fast interrupts)” on page 47 (*processor*)

Additional details about EBB are available here:

<https://github.com/paflib/paflib/wiki/Event-Based-Branching---Overview,-ABI,-and-API>

6.3 Linux operating system-specific optimizations

Next, we describe optimization methods specific to Linux.

6.3.1 GCC, toolchain, and IBM Advance Toolchain

This section describes 32-bit and 64-bit modes and CPU-tuned libraries.

Linux support for 32-bit and 64-bit modes

The compiler and runtime are fully capable of supporting either 32-bit or 64-bit mode applications simultaneously. The compilers can select the target mode through the `-m32` or `-m64` compiler options.

For the SLES and RHEL distributions, the shared libraries have both 32-bit and 64-bit versions. The toolchain (compiler, assembler, linker, and dynamic linker) selects the correct libraries based on the `-m32` or `-m64` option or the mode of the application program.

The Advance Toolchain defaults to 64-bit, as do SLES 11 and RHEL 6. Older distribution compilers defaulted to 32-bit.

Applications can use 32-bit and 64-bit execution modes, depending on their specific requirements, if their dependent libraries are available for the wanted mode.

The 32-bit mode is lighter with a simpler function call sequence and smaller footprint for stack and C++ objects, which can be important for some dynamic language interpreters and applications with many small functions.

The 64-bit mode has a larger footprint because of the larger pointer and general register size, which can be an asset when you handle large data structures or text data, where larger (64-bit) general registers are used for high bandwidth in the memory and string functions.

IBM PowerLinux also supports 64-bit direct memory access (DMA), which can have a significant impact on I/O performance. For details, see *Taking Advantage of 64-bit DMA capability on PowerLinux*, available here:

https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/W51a7ffcf4d4d4b40_9d82_446ebc23c550/page/Taking%20Advantage%20of%2064-bit%20DMA%20capability%20on%20PowerLinux

The handling of floating point and vector data is the same (registers size and format and instructions) for 32-bit and 64-bit modes. Therefore, for these applications, the key decision depends on the address space requirements. For 32-bit Power applications (32-bit mode applications that are running on 64-bit Power hardware with a 64-bit kernel), the address space is limited to 4 GB, which is the limit of a 32-bit address. 64-bit applications are currently limited to 16 TB of application program or data per process. This limitation is not a hardware one, but is a restriction of the shared Linux virtual memory manager implementation. For applications with low latency response requirements, using the larger, 64-bit addressing to avoid I/O latencies using memory mapped files or large local caches is a good trade-off.

CPU-tuned libraries

If an application must support only one Power hardware platform (such as POWER7 and newer), then compiling the entire application with the appropriate `-mcpu=` and `-mtune=` compiler flags might be the best option.

For example, `-mcpu=power7` allows the compiler to use all the POWER7 instructions, such as the Vector Scalar Extended category. The `-mcpu=power7` option also implies `-mtune=power7` if it is not explicitly set.

The GCC compiler does not have any specific POWER7+ optimizations, so just use `-mcpu=power7` or `-mtune=power7`.

The **-mcpu=power8** option allows the compiler to use instructions that were added for the POWER8 processor, such as cryptography built-in functions, direct move instructions that allow data movement between the general purpose registers and the floating point or floating vector registers, and additional vector scalar instructions that were introduced.

-mcpu generates code for a specific machine. If you specify **-mcpu=power7**, the code also runs on a POWER8 machine, but not on a POWER6 machine. Note that **-mcpu=power6x** generates instructions that are not implemented on POWER7 or POWER8 machines, while **-mcpu=power6** generates code that runs on POWER7 and POWER8 machines. The **-mtune** option focuses on optimizing the order of the instructions.

Most applications do need to run on more than one platform, for example, in POWER7 mode and POWER8 mode. For applications composed of a main program and a set of shared libraries or applications that spend significant execution time in other (from the Linux run time or extra package) shared libraries, you can create packages that automatically select the best optimization for each platform.

Linux also supports automatic CPU tuned library selection. There are a number of implementation options for CPU tuned library implementers as described here. For more information, see *Optimized Libraries*, available here:

https://www.ibm.com/developerworks/community/wikis/home?lang=en#/wiki/W51a7ffcf4dfd_4b40_9d82_446ebc23c550/page/Optimized%20Libraries

The Linux Technology Center works with the SUSE and Red Hat Linux Distribution Partners to provide some automatic CPU-tuned libraries for the C/POSIX runtime libraries. However, these libraries might not be supported for all platforms or have the latest optimization.

One advantage of the Advance Toolchain is that the runtime RPMs for the current release do include CPU-tuned libraries for all the currently supported POWER processors and the latest processor-specific optimization and capabilities, which are constantly updated. Additional libraries are added as they are identified. The Advance Toolchain run time can be used with either Advance Toolchain GCC or XL compilers and includes configuration files to simplify linking XL compiled programs with the Advance Toolchain runtime libraries.

These techniques are not restricted to systems libraries, and can be easily applied to application shared library components. The dynamic code path and processor tuned libraries are good starting points. With this method, the compiler and dynamic linker do most of the work. You need only some additional build time and extra media for the multiple library images.

In this example, the following conditions apply:

- ▶ Your product is implemented in your own shared library, such as `libmyapp.so`.
- ▶ You want to support Linux running on POWER5, POWER6, POWER7, and POWER8 systems.
- ▶ DFP and Vector considerations:
 - Your oldest supported platform is POWER5, which does not have a DFP or the Vector unit.
 - POWER6 has DFP and a Vector Unit implementing the older Vector Multimedia eXtension (VMX) (vector float but no vector double) instructions.
 - POWER7 and POWER8 have DFP and the new VSX (the original VMX instructions plus Vector Double and more).
 - Your application benefits greatly from both Hardware Decimal and high performance vector, but if you compile your application with **-mcpu=power7 -O3**, it does not run on POWER5 (no hardware DFP instructions) or POWER6 (no vector double instructions) machines.

You can optimize all three Power platforms if you build and install your application and libraries correctly by completing the following steps:

1. Build the main application binary file and the default version of `libmyapp.so` for the oldest supported platform (in this case, use `-mcpu=power5 -O3`). You can still use decimal data because the Advance Toolchain and the newest SLES 11 and RHEL 6 include a DFP emulation library and run time.
2. Install the application (`myapp`) into the appropriate `./bin` directory and `libmyapp.so` into the appropriate `./lib64` directory. The following paths provide the application main and default run time for your product:

- `/opt/ibm/myapp1.0/bin/myapp`
- `/opt/ibm/myapp1.0/lib64/libmyapp.so`

3. Compile and link `libmyapp.so` with `-mcpu=power6 -O3`, which enables the compiler to generate DFP and VMX instructions for POWER6 machines.
4. Install this version of `libmyapp.so` into the appropriate `./lib64/power6` directory. Here is an example:

```
/opt/ibm/myapp1.0/lib64/power6/libmyapp.so
```

5. Compile and link the fully optimized version of `libmyapp.so` for POWER7 with `-mcpu=power7 -O3`, which enables the compiler to generate DFP and all the VSX instructions. Install this version of `libmyapp.so` into the appropriate `./lib64/power7` directory. Here is an example:

```
/opt/ibm/myapp1.0/lib64/power7/libmyapp.so
```

6. Compile and link the fully optimized version of `libmyapp.so` for POWER8 with `-mcpu=power8 -O3`, which enables the compiler to generate DFP and all the VSX instructions. Install this version of `libmyapp.so` into the appropriate `./lib64/power8` directory. Here is an example:

```
/opt/ibm/myapp1.0/lib64/power8/libmyapp.so
```

By simply running some extra builds, your `myapp1.0` is fully optimized for the current and N-1/N-2 Power hardware releases. When you start your application with the appropriate `LD_LIBRARY_PATH` (including `/opt/ibm/myapp1.0/lib64`), the dynamic linker automatically searches the subdirectories under the library path for names that match the current platform (POWER5, POWER6, POWER7 or POWER8). If the dynamic linker finds the shared library in the subdirectory with the matching platform name, it loads that version; otherwise, the dynamic linker looks in the base `lib64` directory and use the default implementation. This process continues for all directories in the library path and recursively for any dependent libraries.

Using the Advance Toolchain

The latest Advance Toolchain compilers and run time can be downloaded here:

<ftp://ftp.unicamp.br/pub/linuxpatch/toolchain/at>

The latest Advance Toolchain releases (starting with Advance Toolchain 5.0) add multi-core runtime libraries to enable you to take advantage of application level multi-cores. The toolchain currently includes a Power port of the open source version of Intel Thread Building Blocks, the Concurrent Building Blocks software transactional memory library, and the UserRCU library (the application level version of the Linux kernel's Read-Copy-Update concurrent programming technique). Additional libraries are added to the Advance Toolchain run time as needed and if resources allow it.

Linux on Power Enterprise distributions default to 64 KB pages, so most applications automatically benefit from large pages. Larger (16 MB) segments can be best used with the libhugetlbfs API, which is provided with Advance Toolchain. Large segments can be used to back shared memory, malloc storage, and (main) program text and data segments (incorporating large pages for shared library text or data is not supported currently).

6.3.2 Tuning and optimizing malloc

Methods for tuning and optimizing malloc are described in this section.

Linux malloc

Generally, tuning malloc invocations on Linux systems is an application-specific focus.

Improving malloc performance

Linux is flexible regarding the system and application tuning of malloc usage.

By default, Linux manages malloc memory to balance the ability to reuse the memory pool against the range of default sizes of memory allocation requests. Small chunks of memory are managed on the **sbrk** heap. This **sbrk** heap is labeled as [heap] in /proc/self/maps.

When you work with Linux memory allocation, there are a number of tunables available to users. These tunables are coded and used in the Linux malloc.c program. Our examples (“Malloc environment variables” on page 127 and “Linux malloc considerations” on page 128) show two of the key tunables, which force the large sized memory allocations away from using mmap, to using the memory on the program stack by using the **sbrk** system directive.

When you control memory for applications, the Linux operating system automatically makes a choice between using the stack for mallocs with the **sbrk** command, or mmap regions. Mmap regions are typically used for larger memory chunks. When you use mmap for large mallocs, the kernel must zero the newly mmapmed chunk of memory.

Malloc environment variables

Users can define environment variables to control the tunables for a program. The environment variables that are shown in the following examples caused a significant performance improvement across several real-life workloads.

To disable the usage of mmap for mallocs (which includes Fortran allocates), set the max value to zero:

```
MALLOC_MMAP_MAX=0
```

To disable the trim threshold, set the value to negative one:

```
MALLOC_TRIM_THRESHOLD=-1
```

Trimming and using mmap are two different ways of releasing unused memory back to the system. When used together, they change the normal behavior of malloc across C and Fortran programs, which in some cases can change the performance characteristics of the program. You can run one of the following commands to use both actions:

- ▶ # ./my_program
- ▶ # MALLOC_MMAP_MAX=0 MALLOC_TRIM_THRESHOLD=-1 ./my_program

Depending on your application's behavior regarding memory and data locality, this change might do nothing, or might result in performance improvement.

Linux malloc considerations

The Linux GNU C run time includes a default malloc implementation that is optimized for multi-threading and medium sized allocations. For smaller allocations (less than the `MMAP_THRESHOLD`), the default malloc implementation allocates blocks of storage with `sbrk()` called arenas, which are then suballocated for smaller malloc requests. Larger allocations (greater than `MMAP_THRESHOLD`) are allocated by an anonymous mmap, one per request.

The default values are listed here:

<code>DEFAULT_MXFAST</code>	64 (for 32-bit) or 128 (for 64-bit)
<code>DEFAULT_TRIM_THRESHOLD</code>	128 * 1024
<code>DEFAULT_TOP_PAD</code>	0
<code>DEFAULT_MMAP_THRESHOLD</code>	128 * 1024
<code>DEFAULT_MMAP_MAX</code>	65536

Storage within arenas can be reused without kernel intervention. The default malloc implementation uses trylock techniques to detect contentions between POSIX threads, and then tries to assign each thread its own arena. This action works well when the same thread frees storage that it allocates, but it does result in more contention when malloc storage is passed between producer and consumer threads. The default malloc implementation also tries to use atomic operations and more granular and critical sections (lock and unlock) to enhance parallel thread execution, which is a trade-off for better multi-thread execution at the expense of a longer malloc path length with multiple atomic operations per call.

Large allocations (greater than `MMAP_THRESHOLD`) require a kernel syscall for each `malloc()` and `free()`. The Linux Virtual Memory Management (VMM) policy does not allocate any real memory pages to an anonymous `mmap()` until the application touches those pages. The benefit of this policy is that real memory is not allocated until it is needed. The downside is that, as the application begins to populate the new allocation with data, the application experiences multiple page faults, on first touch to allocate and zero fill the page. This situation means that on the initial touching of memory, there is more processing then, as opposed to the earlier timing when the original mmap is done. In addition, this first touch timing can impact the NUMA placement of each memory page.

Such storage is unmapped by `free()`, so each new large malloc allocation starts with a flurry of page faults. This situation is partially mitigated by the larger (64 KB) default page size of the RHEL and SLES on Power Systems; there are fewer page faults than with 4 KB pages.

Malloc tuning parameters

The default malloc implementation provides a `mallopt()` API to allow applications to adjust some tuning parameters. For some applications, it might be useful to adjust the `MMAP_THRESHOLD`, `TOP_PAD`, and `MMAP_MAX` limits. Increasing `MMAP_THRESHOLD` so that most (application) allocations fall below that threshold reduces syscall and page fault impact, and improves application start time. However, this situation can increase fragmentation within the arenas and `sbrk()` storage. Fragmentation can be mitigated to some extent by also increasing `TOP_PAD`, which is the extra memory that is allocated for each `sbrk()`.

Reducing `MMAP_MAX`, which is the maximum number of chunks to allocate with `mmap()`, can also limit the use of `mmap()` when `MMAP_MAX` is set to 0. Reducing `MMAP_MAX` does not always solve the problem. The run time reverts to `mmap()` allocations if `sbrk()` storage, which is the gap between the end of program static data (bss) and the first shared library, is exhausted.

Linux malloc and memory tools

There are several readily available tools in the Linux open source community:

- A website that describes the heap profiler that is used at Google to explore how C++ programs manage memory, found here:

<http://gperftools.googlecode.com/svn/trunk/doc/heapprofile.html>

- *Massif*: a heap profiler, available here:

<http://valgrind.org/docs/manual/ms-manual.html>

For more details about memory management tools, see “Empirical performance analysis using the IBM software development kit (SDK) for PowerLinux” on page 216.

For more information about tuning malloc parameters, see *Malloc Tunable Parameters*, available here:

<http://www.gnu.org/software/libtool/manual/libc/Malloc-Tunable-Parameters.html>

Thread-caching malloc (TCMalloc)

Under some circumstances, an alternative malloc implementation can prove beneficial for improving application performance. Packaged as part of Google's Perftools package (<http://code.google.com/p/gperftools/?redir=1>), and in the Advance Toolchain 5.0.4 release, this specialized malloc implementation can improve performance across a number of C and C++ applications.

TCMalloc uses a thread-local cache for each thread and moves objects from the memory heap into the local cache as needed. Small objects with less than 32 KB are mapped into allocatable size-classes. A thread cache contains a singly linked list of free objects per size-class. Large objects are rounded up to a page size (4 KB) and handled by a central page heap, which is an array of linked lists.

For more information about how TCMalloc works, see *TCMalloc: Thread-Caching Malloc*, available here:

<http://gperftools.googlecode.com/svn/trunk/doc/tcmalloc.html>

The TCMalloc implementation is part of the gperftools project. For more information about this topic, go to this website:

<http://code.google.com/p/gperftools/>

Usage

To use TCMalloc, link TCMalloc into your application by using the `-ltcmalloc` linker flag by running the following command:

```
$ gcc [...] -ltcmalloc
```

You can also use TCMalloc in applications that you did not compile yourself by using `LD_PRELOAD` as follows:

```
$ LD_PRELOAD="/usr/lib/libtcmalloc.so"
```

These examples assume that the TCMalloc library is in `/usr/lib`. With the Advance Toolchain 5.0.4, the 32-bit and 64-bit libraries are in `/opt/at5.0/lib` and `/opt/at5.0/lib64`.

Using TCMalloc with hugepages

To use large pages with TCMalloc, complete the following steps:

1. Set the environment variables for `libhugetlbfs`.
2. Allocate the number of large pages from the system.
3. Set up the `libhugetlbfs` mount point.
4. Monitor large pages usage.

TCMalloc backs up the heap allocation on the large pages only.

Here is a more detailed version of these steps:

1. Set the environment variables for `libhugetlbfs` by running the following commands:

```
- # export TCMALLOC_MEMFS_MALLOCC_PATH=/libhugetlbfs/  
- # export HUGETLB_ELFMAP=RW  
- # export HUGETLB_MORECORE=yes
```

Where:

- `TCMALLOC_MEMFS_MALLOCC_PATH=/libhugetlbfs/` defines the `libhugetlbfs` mount point.
- `HUGETLB_ELFMAP=RW` allocates both RSS and BSS (text/code and data) segments on the large pages, which is useful for codes that have large static arrays, such as Fortran programs.
- `HUGETLB_MORECORE=yes` makes heap usage on the large pages.

2. Allocate the number of large pages from the system by running one of the following commands:

```
- # echo N > /proc/sys/vm/nr_hugepages  
- # echo N > /proc/sys/vm/nr_overcommit_hugepages
```

Where:

- `N` is the number of large pages to be reserved. A peak usage of 4 GB by your program requires 256 large pages (4096/16).
- `nr_hugepages` is the static pool. The kernel reserves $N * 16$ MB of memory from the static pool to be used exclusively by the large pages allocation.
- `nr_overcommit_hugepages` is the dynamic pool. The kernel sets a maximum usage of N large pages and dynamically allocates or deallocates these large pages.

3. Set up the `libhugetlbfs` mount point by running the following commands:

```
- # mkdir -p /libhugetlbfs  
- # mount -t hugetlbfs hugetlbfs /libhugetlbfs
```

4. Monitor large pages usage by running the following command:

```
# cat /proc/meminfo | grep Huge
```

This command produces the following output:

```
HugePages_Total:  
HugePages_Free:  
HugePages_Rsvd:  
HugePages_Surp:  
Hugepagesize:
```

Here we explain these parameters:

- HugePages_Total is the total pages that are allocated on the system for LP usage.
- HugePages_Free is the total free memory available.
- HugePages_Rsvd is the total of large pages that are reserved but not used.
- Hugepagesize is the size of a single LP.

You can monitor large pages by NUMA nodes by running the following command:

```
# watch -d grep Huge /sys/devices/system/node/node*/meminfo
```

MicroQuill SmartHeap

MicroQuill SmartHeap is an optimized malloc that is used for SPECcpu2006 publishes for optimizing performance on selected benchmark components. For more information, see *SmartHeap for SMP: Does your app not scale because of heap contention?*, available here:

<http://www.microquill.com/smartheapsmp/index.html>

6.3.3 Large TOC -mmodel=medium optimization

The Linux Application Binary Interface (ABI) on the Power Architecture is enhanced to optimize larger programs. This ABI both simplifies an application build and improves overall performance.

Previously, the TOC (**-mfull-toc**) defaulted to a single instruction access form that restricts the total size of the TOC to 64 KB. This configuration can cause large programs to fail at compile or link time. Previously, the only effective workaround was to compile with the **-mmimal-toc** option (which provides a private TOC for each source file). The minimal TOC strategy adds a level of indirection that can adversely impact performance.

The **-mmodel=medium** option extends the range of the TOC addressing to +/-2 GB. This setup eliminates most TOC-related build issues. Also, as the Linux ABI TOC includes Global Offset Table (GOT) and local data, you can enable a number of compiler- and linker-based optimizations, including TOC pointer relative addressing for local static and constant data. This setup eliminates a level of indirection and improves the performance of large programs.

Currently, this optimization is available on SLES 11 and RHEL 6 when you are using the system compilers if you use the **-mmodel=medium** option (it is not on by default with those compilers). This optimization is on by default when using Advance Toolchain 4.0 and later.

The medium and large code models are 64-bit only. If you use medium or large code models, you must not use the **-mmimal-toc** option.

6.3.4 POWER7 based distro considerations

For distros that do not recognize the POWER8 processor core, the processor appears as a POWER7 processor, and normal SMT=4 rules apply. This includes the RHEL 6.5 and SLES 11 SP3 releases.

Applications running in this mode still benefit from the efficiencies of the newer processor core, improved cache and memory characteristics, and I/O performance improvements.

6.3.5 Split-core considerations

For POWER8 split-core systems, the CPU numbering can be different. Details are being worked, both for PowerVM based LPARs and KVM Guests.

6.3.6 KVM on Power considerations

The KVM hypervisor host is enabled for SMT=8 support.

A normal POWER8 enabled KVM guest supports SMT=8.

KVM guests can be further limited to SMT=4,2,1 depending on the guest definition (configuration) of the run-time POWER8 mode for that guest.

If the KVM guest is a POWER8-enabled guest which is limited to SMT=4,2,1..., the numbering appears to be consistent with a POWER8 processor. In this case, this means that the `ppc64_cpu` command is unable to turn on more threads.

Information about the topic of operating system-specific optimizations, from the AIX and IBM i perspectives, is available here:

- ▶ 4.3, “AIX operating system-specific optimizations” on page 91 (*AIX*)
- ▶ 5.3, “IBM i operating system-specific optimizations” on page 110 (*IBM i*)

6.4 Related publications

The publications that are listed in this section are considered suitable for a more detailed discussion of the topics that are covered in this chapter:

- ▶ *Red Hat Enterprise Linux 6 Performance Tuning Guide, Optimizing subsystem throughput in Red Hat Enterprise Linux 6, Edition 4.0*, found here:
http://docs.redhat.com/docs/en-US/Red_Hat_Enterprise_Linux/6/html-single/Performance_Tuning_Guide/index.html
- ▶ *SMT settings*, found here:
<http://publib.boulder.ibm.com/infocenter/powersys/v3r1m5/index.jsp?topic=%2Fp7h3%2Fiphc3attributes.htm&resultof%3D%2522%2570%256f%2577%2565%2572%2537%2522%2520%2522%256d%2575%256c%2574%2569%252d%2574%2568%2572%2565%2561%2564%2522%2520>
- ▶ *Simultaneous multithreading*, found here:
<http://pic.dhe.ibm.com/infocenter/lxinfo/v3r0m0/index.jsp?topic=%2Fliiai.hpctune%2Fsmtsetting.htm>
- ▶ *SUSE Linux Enterprise Server System Analysis and Tuning Guide (Version 11 SP3)*, found here:
http://www.suse.com/documentation/sles11/pdfdoc/book_sle_tuning/book_sle_tuning.pdf
- ▶ *POWER6 Decimal Floating Point (DFP)*, found here:
<https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/Power+Systems/page/POWER6+Decimal+Floating+Point+%28DFP%29>
- ▶ *Getting started with OProfile on PowerLinux* (contains references to the `operf` command):
<http://publib.boulder.ibm.com/infocenter/lxinfo/v3r0m0/topic/liacf/oprofgetstart.htm>

- *Power ISA Version 2.07*, found here:

<https://www.power.org/documentation/power-isa-version-2-07/>

See the following sections:

- Section 3.1: Program Priority Registers
- Section 3.2: “or” Instruction
- Section 4.3.4: Program Priority Register
- Section 4.4.3: OR Instruction
- Section 5.3.4: Program Priority Register
- Section 5.4.2: OR Instruction
- Book I – 4 Floating Point Facility
- Book I – 5 Decimal Floating Point
- Book I – 6 Vector Facility
- Book I – 7 Vector-Scalar Floating Point Operations (VSX)
- Book I – Chapter 5 Decimal Floating-Point.
- Book II – 4.2 Data Stream Control Register
- Book II – 4.3.2 Data Cache Instructions
- Book II – 4.4 Synchronization Instructions
- Book II – A.2 Load and Reserve Mnemonics
- Book II – A.3 Synchronize Mnemonics
- Book II – Appendix B. Programming Examples for Sharing Storage
- Book III – 5.7 Storage Addressing



Compilers and optimization tools for C, C++, and Fortran

This chapter describes the optimization and tuning of the POWER8 processor-based server using compilers and tools. It covers the following topics:

- ▶ 7.1, “Compiler versions and optimization levels” on page 136
- ▶ 7.2, “Advanced compiler optimization techniques” on page 137
- ▶ 7.3, “Capitalizing on POWER8 features with the XL and GCC compilers” on page 142
- ▶ 7.4, “IBM Feedback Directed Program Restructuring (FDPR)” on page 149
- ▶ 7.5, “Using the Advance Toolchain with IBM XLC and XLF” on page 158
- ▶ 7.6, “Related publications” on page 159

7.1 Compiler versions and optimization levels

The IBM XL compilers are updated periodically to improve application performance and add processor-specific tuning and capabilities. The XLC13 and XLF15 compilers for AIX and Linux are the first versions to include the capabilities of POWER8, and are the preferred version for projects that target current generation systems.

For the GNU GCC, G++ and gfortran compilers on Linux, the IBM Advance Toolchain 7.0 (GCC-4.8) has POWER8 enabled. Note that the normal distro version of GCC-4.8 does not have POWER8 support. It is anticipated that a future distro GCC release will contain POWER8 support. XLF is preferred over gfortran for its high floating point performance characteristics.

For all production codes, it is imperative to enable a minimum level of compiler optimization by adding the **-O** option for the XL compilers, or **-O2** with the GNU compilers (**-O3** is the preferred option). Without optimization, the focus of the compiler is on faster compilation and debug ability, and it generates code that performs poorly at run time. In practice, many projects set up a dual build environment, with a development build without optimization for use during development and debugging, and a production build with optimization to be used for performance verification and production delivery.

For projects with increased focus on runtime performance, take advantage of the more advanced compiler optimization. For numerical or compute-intensive codes, the XL compiler options **-O3** or **-qhot -O3** enable loop transformations, which improve program performance by restructuring loops to make their execution more efficient by the target system. These options perform aggressive transformations that can sometimes cause minor differences on precision of floating point computations. If that is a concern, the original program semantics can be fully recovered with the **-qstrict** option.

For GCC, the minimum suggested level of optimization is **-O3**. The GCC default is a strict mode, but the **-ffast-math** option disables strict mode. The **-Ofast** option combines **-O3** with **-ffast-math** in a single option. Other important options include **-fpeel-loops**, **-funroll-loops**, **-ftree-vectorize**, **-fvect-cost-model**, and **-mmodel=medium**.

By default, these compilers generate code that run on various Power Systems. Options should be added to exclude older processor chips that are not supported by the target application. This configuration might enable better code generation as the compiler takes advantage of capabilities not available on those older systems.

There are two major XL compiler options to control this support:

- ▶ **-qarch**: Indicates the oldest processor chip generation that the binary file supports.
- ▶ **-qtune**: Indicates the processor chip generation of most interest for performance.

For example, for an application that must run on POWER7 systems, but for which most users are on a POWER8 system, the appropriate combination is **-qarch=pwr7 -qtune=pwr8**. For an application that must run well across both POWER7 and POWER8 systems in current common usage, consider using **-qtune=balanced**.

On GCC, the equivalent options are **-mcpu** and **-mtune**. So, for an application that must run on POWER7, but which is usually run on POWER8, the options are **-mcpu=power7** and **-mtune=power8**.

The XLC13 and XLF15 compilers for AIX and Linux introduce an extension to `-qtune` to indicate the SMT mode that the application will most often run in. For example, for an application that must run on POWER8 systems and will most often run in SMT4 mode (for example, four hardware threads per core), use `-qarch=pwr8 -qtune=pwr8:smt4`. If the same application might run in several different SMT modes, consider using `-qarch=pwr8 -qtune=pwr8:balanced`.

The POWER8 processor supports the Vector Scalar eXtension (VSX) instruction set, which improves performance for numerical applications over regular data sets. These performance features can increase the performance of some computations, and can be accessed manually by using the AltiVec vector extensions, or automatically by the XL compiler by using `-O3 -qhot` options with `-qarch=pwr7` or `-qarch=pwr8`. By default, these options implicitly enable `-qsimd` which allows the XL compilers to transform loops in an application to exploit VSX instructions. The POWER8 processor includes several extensions to the Vector Multimedia eXtension (VMX) and VSX instruction sets which can improve performance of applications using 64-bit integer types and single-precision floating point.

The GCC compiler equivalents are the `-maltivec` and `-mvsx` options, which you should combine with `-ftree-vectorize` and `-fvect-cost-model`. On GCC, the combination of `-O3` and `-mcpu=power7` or `-mcpu=power8` implicitly enables AltiVec and VSX code generation with auto-vector (`-ftree-vectorize`) and `-mpopcntd`. Other important options include `-mrecip=rsqrt` and `-mveclibabi=mass` (which *require* `-ffast-math` or `-Ofast` to be effective). If the compiler uses optimizations dependent on the MASS libraries, the link command must explicitly name the MASS library directories and library names.

For more information about this topic, see 7.6, “Related publications” on page 159.

7.2 Advanced compiler optimization techniques

This section describes some of the more advanced compiler optimization techniques.

7.2.1 Common prerequisites

Compiler analysis and transformations improve runtime performance by changing the translation of the program source into assembly code. Changes in these translations might cause the application to behave differently, possibly even causing it to produce incorrect results.

Compilers follow rules and assumptions that are part of the programming language to perform this transformation. If the programmer breaks some of these rules, it is possible for the application to misbehave, and it might do so only at higher optimization levels, where it is more difficult for the problem to be diagnosed.

To put this situation into perspective, imagine a C program with three variables: “`int a[4], b, c;`”. These variables are normally placed contiguously in memory. If the user runs a statement of the form `a[5]=0`, this statement breaks the language rules, but if variable `b` is unused, the statement might overwrite variable `b` and the program might continue to behave correctly. However, if, at a higher optimization level, variable `b` is eliminated, as the compiler determines it is unused, the incorrect statement might overwrite variable `c`, triggering a runtime failure.

It is critical, then, to eliminate programming errors as higher optimization is applied. Testing the application thoroughly without optimization is a good initial step, but it is not required or sufficient. The application must be tested at the optimization level to be used in production.

7.2.2 XL compiler family

There are several XL compiler programming errors that can provide guidance toward optimization.

Prerequisites

The XL compilers assist with identifying certain programming errors that are outlined 7.2.1, “Common prerequisites” on page 137:

- ▶ Static analysis/warnings: The XL compilers can identify suspicious code constructs, and provide some information about these constructs through the **-qinfo=all** option. Examine the output of this option to identify suspicious code constructs and validate that the constructs are correct.
- ▶ Runtime analysis or warning: The XL compilers can cause the application to perform runtime checks to validate program correctness by using the **-qcheck** option. This option triggers a program abort when an error condition (such as a null pointer dereference or out-of-bounds array access) is run, identifying a problem and making it easier for you to identify it. This option has a significant performance cost, so use it only during functional verification, not on a production environment.
- ▶ Aliasing compliance: The C, C++, and Fortran languages specify rules that govern the access of data through overlapping pointers. These rules are brought into play aggressively by optimization techniques, but they can lead to incorrect results if they are broken. The compiler can be instructed not to take advantage of these rules, at a cost of runtime performance. This situation can be useful for older code that is written without following these rules. The options to request this optimization are **-qalias=noansi** for C/C++ and **-qalias=nostd** for Fortran.

The XLC13 and XLF15 compilers include enhancements to **-qinfo** and **-qcheck** to detect accesses to uninitialized variables and stack corruption or stack clobbering.

High-order transformations (HOT)

The XL compilers have sophisticated optimizations to improve the performance of numeric applications. These applications often contain regular loops that process large amounts of data. The high-order transformations (HOT) optimizations in these compilers analyze these loops, identify opportunities for restructuring them to improve cache usage, improve data reuse, and expose more instruction-level parallelism to the hardware. For these types of applications, the performance impact of this option can be substantial.

There are two levels of aggressiveness to the HOT optimization framework in these compilers:

- ▶ Level 0, which is the default at optimization level **-O3**, performs a minimal amount of loop optimization, focusing on simple opportunities while it minimizes compilation time.
- ▶ Level 1, which is the default at optimization levels **-O4** and up, performs full loop analysis and transformation of loops.

The HOT optimizations can be explicitly requested through the **-qhot=level=0** and **-qhot=level=1** options. The **-qhot** option alone enables **-qhot=level=1**. The **-O3 -qhot** options are preferred for numerical applications.

OpenMP

The OpenMP API is an industry specification for shared-memory parallel programming. The latest XL Compilers provide a full implementation of the OpenMP 3.0 specification in C, C++, and Fortran. You can program with OpenMP to capitalize on the incremental introduction of parallelism in an existing application by adding pragmas or directives to specify how the application can be parallelized.

For applications with available parallelism, OpenMP can provide a simple solution for parallel programming, without requiring low-level thread manipulation. The OpenMP implementation on the XL compilers is available by using the **-qsmp=omp** option.

Whole-program analysis (IPA)

Traditional compiler optimizations operate independently on each application source file. Inter-procedural optimizations operate at the whole-program scope, using the interaction between parts of the application on different source files. It is often effective for large-scale applications that are composed of hundreds or thousands of source files.

On the XL compilers, these capabilities are accessed by using the **-qipa** option. It is also implied when you use optimization levels **-O4** and **-O5**. In this phase, the compiler saves a high-level representation of the program in the object files during compilation, and reoptimizes it at the whole-program scope during the link phase. For this situation to occur, the compiler driver must be used to link the resulting binary, instead of invoking the system linker directly.

Whole-program analysis (IPA) is effective on programs that use many global variables, overflowing the default AIX limit on global symbols. If the application requires the use of the **-bbigtoc** option to link successfully on AIX, it is likely a good candidate for IPA optimization.

There are three levels of IPA optimization on the XL compilers (0, 1, and 2). By default, **-qipa** implies **ipa=level=1**, which performs basic program restructuring. For more aggressive optimization, apply **-qipa=level=2**, which performs full program restructuring during the link step. The time that it takes to complete the link step can increase significantly.

Optimization that is based on Profile Directed Feedback

Profile-based optimization allows the compiler to collect information about the program behavior and use that information when you make code generation decisions. It involves compiling the program twice: first, to generate an *instrumented* version of the application that collects program behavior data when run, and a second time to generate an optimized binary file using information that is collected by running the instrumented binary through a set of typical inputs for the application.

Profile-based optimization in the XL compiler is accessed through the **-qpdf1** and **-qpdf2** options, on top of **-O** or higher optimization levels. The instrumented binary file is generated by using **-qpdf1** on top of all other options, and the resulting binary file generates the profile data on a file, named **._pdf** by default.

The Profile Directed Feedback (PDF) framework on the XL compilers is built on top of the IPA infrastructure, with **-qpdf1** and **-qpdf2** implying **-qipa=level=0**. For the PDF2 step, it is possible to reuse the object files from the **-qpdf1** compilation step, and relink only the application with the **-qpdf2** option.

For PDF optimizations to be successful, the instrumented workload must be run with common workloads that reflect common usage of the application. Use multiple workloads that can exercise the program in different ways. The data for all instrumentation runs are aggregated into a single PDF file and used during optimization.

For the PDF profile data to be written out at the end of execution, the program must either implicitly or explicitly call the **exit()** library subroutine. Using **exit()** causes code that is introduced as part of the PDF instrumentation to be run and write out the PDF profile data. In contrast, running the **_exit()** system call skips the writing of the PDF profile data file, which results in inaccurate profile data being recorded.

7.2.3 GCC compiler family

The information in this section applies specifically to the GCC compiler family.

Prerequisites

The GCC compiler assists with identifying certain programming errors that are outlined in 7.2.1, “Common prerequisites” on page 137:

- ▶ Static analysis and warnings. The **-pedantic** and **-pedantic-errors** options warn of violations of ISO C or ISO C++ standards.
- ▶ The language standard to enforce and the aliasing compliance requirements are specified by the **-std**, **-ansi**, and **-fno-strict-aliasing** options. Here is an example:
 - ISO C 1990 level: **-std=c89**, **-std=iso9899:1990**, and **-ansi**
 - ISO C 1998 level: **-std=c99** and **-std=iso9899:1999**
 - Do not assume strict aliasing rules for the language level: **-fno-strict-aliasing**

The GCC compiler documentation contains more details about these options.^{1, 2, 3}

High order transformations (HOT)

The GCC compilers have sophisticated additional optimizations beyond **-O3** to improve the performance of numeric applications. These applications often contain regular loops that process large amounts of data. These optimizations, when enabled, analyze these loops, identify opportunities for restructuring them to improve cache usage, improve data reuse, and expose more instruction-level parallelism to the hardware. For these types of applications, the performance impact of this option can be substantial. The key compiler options include:

- ▶ **-fpeel-loops**
- ▶ **-funroll-loops**
- ▶ **-ftree-vectorize**
- ▶ **-fvect-cost-model**
- ▶ **-mcmmodel=medium**

Specifying the **-mveclibabi=mass** option and linking to the MASS libraries enables more loops for **-ftree-vectorize**. The MASS libraries support only static archives for linking, and so they require explicit naming and library search order for each platform/mode:

- ▶ POWER8 32-bit: **-L<MASS-dir>/lib -lmassvp8 -lmass_simdp8 -lmass -lm**
- ▶ POWER8 64-bit: **-L<MASS-dir>/lib64 -lmassvp8_64 -lmass_simdp8_64 -lmass_64 -lm**
- ▶ POWER7 32-bit: **-L<MASS-dir>/lib -lmassvp7 -lmass_simdp7 -lmass -lm**
- ▶ POWER7 64-bit: **-L<MASS-dir>/lib64 -lmassvp7_64 -lmass_simdp7_64 -lmass_64 -lm**
- ▶ POWER6 32-bit: **-L<MASS-dir>/lib -lmassvp6 -lmass -lm**
- ▶ POWER6 64-bit: **-L<MASS-dir>/lib64 -lmassvp6_64 -lmass_64 -lm**

ABI improvements

The **-mcmmodel={medium|large}** option implements important ABI improvements that are further optimized in hardware for future generations of the POWER processor. This optimization extends the Table-Of-Content (TOC) to 2 GB and eliminates the previous requirement for **-mmminimal-toc** or multi-TOC switching within a single a program or library. The default for newer GCC compilers (including Advance Toolchain 4.0 and later) is **-mcmmodel=medium**. This model logically extends the TOC to include local static data and constants and allows direct data access relative to the TOC pointer.

¹ *Language Standards Supported by GCC*, available here:

<http://gcc.gnu.org/onlinedocs/gcc-3.4.2/gcc/Standards.html#Standards>

² *Options Controlling C Dialect*, available here:

<http://gcc.gnu.org/onlinedocs/gcc-3.4.2/gcc/C-Dialect-Options.html#C-Dialect-Options>

³ *Options That Control Optimization, and specifically the discussion of -fstrict-aliasing*, available here:

<http://gcc.gnu.org/onlinedocs/gcc-3.4.2/gcc/Optimize-Options.html#Optimize-Options>

OpenMP

The OpenMP API is an industry specification for shared-memory parallel programming. The current GCC compilers, starting with GCC- 4.4 (Advance Toolchain 4.0 and later), provide a full implementation of the OpenMP 3.0 specification in C, C++, and Fortran. Programming with OpenMP allows you to benefit from the incremental introduction of parallelism in an existing application by adding pragmas or directives to specify how the application can be parallelized.

For applications with available parallelism, OpenMP can provide a simple solution for parallel programming, without requiring low-level thread manipulation. The GNU OpenMP implementation on the GCC compilers is available under the **-fopenmp** option. GCC also provides auto-parallelization under the **-ftree-parallelize-loops** option.

Whole-program analysis (IPA)

Traditional compiler optimizations operate independently on each application source file. Inter-procedural optimizations operate at the whole-program scope, using the interaction between parts of the application on different source files. It is often effective for large-scale applications that are composed of hundreds or thousands of source files.

Starting with GCC- 4.6 (Advance Toolchain 5.0), there is the Link Time Optimization (LTO) feature. LTO allows separate compilation of multiple source files but saves additional (abstract program description) information in the resulting object file. Then, at application link time, the linker can collect all the objects (with additional information) and pass them back to the compiler (GCC) for whole program IPA and final code generation.

The GCC LTO feature is enabled on the compile and link phases by the **-flto** option. A simple example follows:

```
gcc -flto -O3 -c a.c
gcc -flto -O3 -c b.c
gcc -flto -o program a.o b.o
```

Here are some additional options that can be used with **-flto**:

- ▶ **-flto-partition={lto1|balanced|none}**
- ▶ **-flto-compression-level=*n***

Detailed descriptions about **-flto** and its related options are in *Options That Control Optimization*, available here:

<http://gcc.gnu.org/onlinedocs/gcc-4.6.3/gcc/Optimize-Options.html#Optimize-Options>

Profiled-based optimization

Profile-based optimization allows the compiler to collect information about the program behavior and use that information when you make code generation decisions. It involves compiling the program twice: first, to generate an *instrumented* version of the application that collects program behavior data when run, and a second time to generate an optimized binary using information that is collected by running the instrumented binary through a set of typical inputs for the application.

Profile-based optimization in the GCC compiler is accessed through the **-fprofile-generate** and **-fprofile-use** options on top of **-O2** optimization levels. The instrumented binary is generated by using **-fprofile-generate** on top of all other options, and the resulting binary file generates the profile data in a file, named `._pdf` by default. Here is an example:

```
gcc -fprofile-generate -O3 -c a.c
gcc -fprofile-generate -O3 -c b.c
gcc -fprofile-generate -o program a.o b.o
```

```

program < sample1
program < sample2
program < sample3
gcc -fprofile-use -O3 -c a.c
gcc -fprofile-use -O3 -c b.c
gcc -fprofile-use -o program a.o b.o

```

These are additional options that are related to GCC PDF:

- fprofile-correction** Corrects for missing counter samples from multi-threaded applications.
- fprofile-dir=PATH** Specifies the directory for generating and using profile data.
- fprofile-generate=PATH** Combines -fprofile-generate and -fprofile-dir.
- fprofile-use=PATH** Combines -fprofile-use and -fprofile-dir.

Detailed descriptions about **-fprofile-generate** and its related options can be found *Options That Control Optimization*, available here:

<http://gcc.gnu.org/onlinedocs/gcc-4.6.3/gcc/Optimize-Options.html#Optimize-Options>

For more information about this topic, see 7.6, “Related publications” on page 159.

7.3 Capitalizing on POWER8 features with the XL and GCC compilers

This section describes built-in functions provided by the XL and GCC compiler families for high-level language access to new POWER8 features and instructions.

7.3.1 In-core cryptography

The GCC, XL C/C++, and XL Fortran compilers provide built-in functions for the in-core cryptography instructions. For GCC, the following built-in functions require **-mcpu=power8** or **-mcrypto**. For the XL compiler family, **-qarch=pwr8** is required.

AES

The following built-in functions are provided for implementation of the AES algorithm:

- **vsbox**
 - GCC: vector unsigned long long **__builtin_crypto_vsbox** (vector unsigned long long)
 - XLC/C++: vector unsigned char **__vsbox** (vector unsigned char) XLF: **VSBOX** (ARG1), where ARG1 and result are unsigned vector types of kind 1
- **vcipher**
 - GCC: vector unsigned long long **__builtin_crypto_vcipher** (vector unsigned long long, vector unsigned long long)
 - XLC/C++: vector unsigned char **__vcipher** (vector unsigned char, vector unsigned char)
 - XLF: **VCIPHER** (ARG1,ARG2), where ARG1, ARG2 and result are unsigned vector types of kind 1

► **vcipherlast**

- GCC: vector unsigned long long **__builtin_crypto_vcipherlast** (vector unsigned long long, vector unsigned long long)
- XLC/C++: vector unsigned char **__vcipherlast** (vector unsigned char, vector unsigned char)
- XLF: **VCIPHERLAST** (ARG1,ARG2), where ARG1, ARG2 and result are unsigned vector types of kind 1

► **vncipher**

- GCC: vector unsigned long long **__builtin_crypto_vncipher** (vector unsigned long long, vector unsigned long long)
- XLC/C++: vector unsigned char **__vncipher** (vector unsigned char, vector unsigned char)
- XLF: **VNCIPHER** (ARG1,ARG2), where ARG1, ARG2 and result are unsigned vector types of kind 1

► **vncipherlast**

- GCC: vector unsigned long long **__builtin_crypto_vncipherlast** (vector unsigned long long, vector unsigned long long)
- XLC/C++: vector unsigned char **__vncipherlast** (vector unsigned char, vector unsigned char)
- XLF: **VNCIPHERLAST** (ARG1,ARG2), where ARG1, ARG2 and result are unsigned vector types of kind 1

See “AES” on page 43 for further information.

AES Galois Counter Mode (GCM)

The following built-in functions are provided for implementation of the Galois Counter Mode of AES:

► **vpmsumd**

- GCC: vector unsigned long long **__builtin_crypto_vpmsum** (vector unsigned long long, vector unsigned long long)
- XLC/C++: vector unsigned long long **__vpmsumd** (vector unsigned long long, vector unsigned long long)
- XLF: **VPMSUMD** (ARG1, ARG2), where ARG1, ARG2 and result are unsigned vector types of kind 8

See “AES special mode of operation: Galois Counter Mode (GCM)” on page 43 for further information.

SHA-2

The following built-in functions are provided for implementation of SHA-2 hash functions:

► **vshasigmad**

- GCC: vector unsigned long long **__builtin_crypto_vshasigmad** (vector unsigned long long, int, int)
- XLC/C++: vector unsigned long long **__vshasigmad** (vector unsigned long long, int, int)
- XLF: **VSHASIGMAD** (ARG1,ARG2,ARG3), where ARG1 and result are unsigned vector types of kind 8, and ARG2, ARG3 are integer types

► **vshasigmaw**

- GCC: vector unsigned int **__builtin_crypto_vshasigmaw** (vector unsigned int, int, int)
- XLC/C++: vector unsigned int **__vshasigmaw** (vector unsigned int, int, int)
- XLF: **VSHASIGMAW** (ARG1, ARG2, ARG3), where ARG1 and result are unsigned vector types of kind 4, and ARG2, ARG3 are integer types

See “SHA-2” on page 43 for further information.

CRC

The following built-in functions are provided for implementation of the CRC algorithm:

► **vpmsumd**

- GCC: vector unsigned long long **__builtin_crypto_vpmsumd** (vector unsigned long long, vector unsigned long long)
- XLC/C++: vector unsigned long long **__vpmsumd** (vector unsigned long long, vector unsigned long long)
- XLF: **VPMSUMD** (ARG1, ARG2), where ARG1, ARG2 and result are unsigned vector types of kind 8

► **vpmsumw**

- GCC: vector unsigned int **__builtin_crypto_vpmsum** (vector unsigned int, vector unsigned int)
- XLC/C++: vector unsigned int **__vpmsumw** (vector unsigned int, vector unsigned int)
- XLF: **VPMSUMW** (ARG1, ARG2), where ARG1, ARG2 and result are unsigned vector types of kind 4

► **vpmsumh**

- GCC: vector unsigned short **__builtin_crypto_vpmsum** (vector unsigned short, vector unsigned short)
- XLC/C++: vector unsigned short **__vpmsumh** (vector unsigned short, vector unsigned short)
- XLF: **VPMSUMH** (ARG1, ARG2), where ARG1, ARG2 and result are unsigned vector types of kind 2

► **vpmsumb**

- GCC: vector unsigned char **__builtin_crypto_vpmsum** (vector unsigned char, vector unsigned char)
- XLC/C++: vector unsigned char **__vpmsumb** (vector unsigned char, vector unsigned char)
- XLF: **VPMSUMB** (ARG1, ARG2), where ARG1, ARG2 and result are unsigned vector types of kind 1

Information about the topic of in-core cryptography, from the processor and OS perspectives, is available here:

- 2.2.7, “In-core cryptography and integrity enhancements” on page 42 (*processor*)
- 4.2.7, “On-chip encryption accelerator” on page 85 (*AIX*)

7.3.2 Compiler support for VSX

XLC supports vector processing technologies through language extensions on both AIX and Linux. GCC supports using the VSX engine on Linux. XL and GCC C implement and extend the AltiVec Programming Interface specification. In the extended syntax, type qualifiers and storage class specifiers can precede the keyword vector (or its alternative spelling, `__vector`) in a declaration.

Also, the XL compilers are able to automatically generate VSX instructions from scalar code when they generate code that targets the POWER7 processor. This task is accomplished by using the `-qsimd=auto` option with the `-O3` optimization level or higher. GCC also automatically generates VSX instructions from scalar code when generating code for POWER7 (or later). This is accomplished by specifying `-O3` (or `-ftree-vectorize`).

Table 7-1 lists the supported vector data types and the size and possible values for each type.

Table 7-1 Vector data types

Type	Interpretation of content	Range of values
vector unsigned char	16 unsigned char	0..255
vector signed char	16 signed char	-128..127
vector bool char	16 unsigned char	0, 255
vector unsigned short	8 unsigned short	0..65535
vector unsigned short int		
vector signed short	8 signed short	-32768..32767
vector signed short int		
vector bool short	8 unsigned short	0, 65535
vector bool short int		
vector unsigned int	4 unsigned int	0..2 ³² -1
vector unsigned long	4 unsigned int (32-bit)	0..2 ³² -1
vector unsigned long int	2 unsigned long int (64-bit)	0..2 ⁶⁴ -1
vector signed int	4 signed int	-2 ³¹ ..2 ³¹ -1
vector signed long	4 signed int (32-bit)	-2 ³¹ ..2 ³¹ -1
vector signed long int	2 signed long int (64-bit)	-2 ⁶³ ..2 ⁶³ -1
vector bool int	4 unsigned int	0, 2 ³² -1
vector bool long	4 unsigned int (32-bit)	0, 2 ³² -1
vector bool long int	2 unsigned long int (64-bit)	0, 2 ⁶⁴ -1
vector float	4 float	IEEE-754 single (32 bit) precision floating point values
vector double	2 double	IEEE-754 double (64 bit) precision floating point values
vector pixel	8 unsigned short	1/5/5/5 pixel

Vector types: The *vector double* type requires architectures that support the VSX instruction set extensions, such as POWER7. You must specify the XL `-qarch=pwr7 -qaltivec` compiler options when you use this type, or the GCC `-mcpu=power7` or `-mvsx` options.

The hardware does not have instructions for supporting vector unsigned long long, vector bool long long, or vector signed long long. In GCC, you can declare these types, but the only hardware operation you can use these types for is vector floating point convert. In 64-bit mode, vector long is the same as vector long long. In 32-bit mode, these types are not permitted.

All vector types are aligned on a 16-byte boundary. An aggregate that contains one or more vector types is aligned on a 16-byte boundary, and padded, if necessary, so that each member of vector type is also 16-byte aligned. Vector data types can use some of the unary, binary, and relational operators that are used with primitive data types. All operators require compatible types as operands unless otherwise stated. For more information about the operator's usage, see the XLC online publications^{4, 5, 6}.

Individual elements of vectors can be accessed by using the VMX or the VSX built-in functions. For more information about the VMX and the VSX built-in functions, see the built-in functions section of *Vector Built-in Functions*.⁷s

Vector initialization

A vector type is initialized by a vector literal or any expression that has the same vector type. Here is an example:⁸

```
vector unsigned int v1;
vector unsigned int v2 = (vector unsigned int)(10); // XL only, not GCC
v1 = v2;
```

The number of values in a braced initializer list must be less than or equal to the number of elements of the vector type. Any uninitialized element is initialized to zero.

Here are examples of vector initialization using initializer lists:

```
vector unsigned int v1 = {1}; // initialize the first 4 bytes of v1 with 1
                             // and the remaining 12 bytes with zeros
vector unsigned int v2 = {1,2}; // initialize the first 8 bytes of v2 with 1 and 2
                             // and the remaining 8 bytes with zeros
vector unsigned int v3 = {1,2,3,4}; // equivalent to the vector literal
                             // (vector unsigned int) (1,2,3,4)
```

⁴ Support for POWER7 processors, available here:

<http://publib.boulder.ibm.com/infocenter/comphelp/v111v131/index.jsp?topic=/com.ibm.xlc111.aix.doc/gtstart/architecture.html>

⁵ Vector built-in functions, available here:

http://publib.boulder.ibm.com/infocenter/comphelp/v111v131/index.jsp?topic=/com.ibm.xlc111.aix.doc/compiler_ref/vec_intrin_cpp.html

⁶ Initialization of vectors (IBM extension), available here:

http://pic.dhe.ibm.com/infocenter/comphelp/v111v131/index.jsp?topic=%2Fcom.ibm.xlcpp111.aix.doc%2Flanguage_ref%2Fvector_init.html

⁷ Vector built-in functions, available here:

http://pic.dhe.ibm.com/infocenter/comphelp/v111v131/index.jsp?topic=%2Fcom.ibm.xlcpp111.aix.doc%2Fcompiler_ref%2Fvec_intrin_cpp.html

⁸ Vector types (IBM extension), available here:

http://pic.dhe.ibm.com/infocenter/comphelp/v111v131/index.jsp?topic=%2Fcom.ibm.xlc111.aix.doc%2Flanguage_ref%2Faltivec_types.html

How to use vector capability in POWER8

When you target a POWER processor that supports VMX or VSX, you can request the compiler to transform code into VMX or VSX instructions. These machine instructions can run up to 16 operations in parallel. This transformation mostly applies to loops that iterate over contiguous array data and perform calculations on each element. You can use the NOSIMD directive to prevent the transformation of a particular loop:⁹

- ▶ Using a compiler: Compiler versions that recognize the POWER8 architecture are XL C/C++ 13.1 and XLF Fortran 15.1 or recent versions of GCC, including the Advance Toolchain, and the SLES 11SP1 or Red Hat RHEL6 GCC compilers:
 - For C:
 - `xlc -qarch=pwr8 -qtune=pwr8 -O3 -qhot`
 - `gcc -mcpu=power8 -mtune=power8 -O3`
 - For Fortran
 - `xlf -qarch=pwr8 -qtune=pwr8 -O3 -qhot`
 - `gfortran -mcpu=power8 -mtune=power8 -O3`
- ▶ Using Engineering and Scientific Subroutine (ESSL) libraries with vectorization support:
 - Select routines have vector analogs in the library
 - Key FFT, BLAS routines

Information about the topic of VSX, from the processor and OS perspectives, is available here:

- ▶ 2.2.5, “Vector Scalar eXtension (VSX)” on page 40 (*processor*)
- ▶ 4.2.5, “Vector Scalar eXtension (VSX)” on page 82 (*AIX*)
- ▶ 5.2.3, “Vector Scalar eXtension (VSX)” on page 103 (*IBM i*)
- ▶ 6.2.5, “Vector Scalar eXtension (VSX)” on page 120 (*Linux*)

7.3.3 Built-in functions for storage synchronization

The XL C/C++ compiler provides built-in functions for direct usage of the storage synchronization load/store operations, as described in 2.2.9, “Storage synchronization (sync, lwsync, lwarx, stwcx, and eieio)” on page 44. New functions for POWER8 are indicated and require `-qarch=pwr8`.

Each pair of builtins is used to implement a read-modify-write operation on a memory location of a given size, where the load built-in (beginning with `__l`) returns the loaded value, and the store built-in (beginning with `__st`) returns 1 if the store succeeds, 0 otherwise. The functions work together to guarantee that if the store succeeds, then no other processor or mechanism can modify the target between the time of the load execution and store completion.

```
char __lbarx(volatile char* addr)
int __stbcx(volatile char* addr, char data)
```

Load and store, respectively, the byte value located at `addr`. Requires `-qarch=pwr8`.

```
short __lharx(volatile short* addr)
int __sthcx(volatile short* addr, short data)
```

Load and store, respectively, the halfword value located at `addr`. `addr` must be aligned on a halfword boundary. Requires `-qarch=pwr8`.

```
int __lwarx(volatile int* addr)
```

⁹ Ibid

```
int __stwcx(volatile int* addr, int data)
```

Load and store, respectively, the word value located at **addr**. **addr** must be aligned on a word boundary.

```
long __ldarx(volatile long* addr)
```

```
long __stdcx(volatile long* addr, long data)
```

Load and store, respectively, the doubleword value located at **addr**. **addr** must be aligned on a doubleword boundary. Only valid in 64-bit mode.

```
void __lqarx(volatile long* addr, long data[2])
```

```
long __stqcx(volatile long* addr, long data[2])
```

__lqarx loads the quadword value at **addr** into the quadword location specified by data. **__stqcx** stores the quadword value in data to the quadword location specific by **addr**. Both **addr** and data must be aligned on a quadword boundary. Only valid in 64-bit mode, and requires **-qarch=pwr8**.

7.3.4 Data Streams Control Register (DSCR) controls

The XL C/C++ and XL Fortran compilers provide the following built-in functions to modify the DSCR setting. The **-qarch=pwr8** option is required to use the following built-in functions. For descriptions of the DSCR fields, see Table 2-4 on page 35.

- ▶ C/C++: **void __software_transient_enable(int)**
- ▶ Fortran: **SOFTWARE_TRANSIENT_ENABLE(flag)**, where flag is a scalar of type logical
 - Set the SWTE bit to the provided value (0 or 1)
- ▶ C/C++: **void __hardware_transient_enable(int)**
- ▶ Fortran: **HARDWARE_TRANSIENT_ENABLE(flag)**, where flag is a scalar of type logical
 - Set the HWTE bit to the provided value (0 or 1)
- ▶ C/C++: **void __store_transient_enable(int)**
- ▶ Fortran: **STORE_TRANSIENT_ENABLE(flag)**, where flag is a scalar of type logical
 - Set the STE bit to the provided value (0 or 1)
- ▶ C/C++: **void __load_transient_enable(int)**
- ▶ Fortran: **LOAD_TRANSIENT_ENABLE(flag)**, where flag is a scalar of type logical
 - set the LTE bit to the provided value (0 or 1)
- ▶ C/C++: **void __software_unit_count_enable(int)**
- ▶ Fortran: **SOFTWARE_UNIT_COUNT_ENABLE(flag)**, where flag is a scalar of type logical
 - Set the SWUE bit to the provided value (0 or 1)
- ▶ C/C++: **void __hardware_unit_count_enable(int)**
- ▶ Fortran: **HARDWARE_UNIT_COUNT_ENABLE(flag)**, where flag is a scalar of type logical
 - Set the HWUE bit to the provided value (0 or 1)
- ▶ C/C++: **void __set_prefetch_unit_count(int)**
- ▶ Fortran: **SET_PREFETCH_UNIT_COUNT(cnt)**, where cnt is a scalar of type integer
 - Set the UNITCNT field to the provided value (in range [0,1023])
- ▶ C/C++: **void __depth_attainment_urgency(int)**

- ▶ Fortran: **DEPTH_ATTAINMENT_URGENCY(cnt)**, where cnt is a scalar of type integer
 - Set the URG field to the provided value (in range [0,7])
- ▶ C/C++: **void __load_stream_disable(int)**
- ▶ Fortran: **LOAD_STREAM_DISABLE(flag)**, where flag is a scalar of type logical
 - Set the LSD bit to the provided value (0 or 1)
- ▶ C/C++: **void __stride_n_stream_enable(int)**
- ▶ Fortran: **STRIDE_N_STREAM_ENABLE(flag)**, where flag is a scalar of type logical
 - Set the SNSE bit to the provided value (0 or 1)
- ▶ C/C++: **void __default_prefetch_depth(int)**
- ▶ Fortran: **DEFAULT_PREFETCH_DEPTH(cnt)**, where cnt is a scalar of type integer
 - Set the DPFDP field to the provided value (in range [0,7])
- ▶ C/C++: **unsigned long long __prefetch_get_dscr_register(void)**
- ▶ Fortran: **PREFETCH_GET_DSCR_REGISTER()**, with return type integer*8
 - Get the current 64-bit DSCR register value
- ▶ C/C++: **void __prefetch_set_dscr_register(void)**
- ▶ Fortran: **PREFETCH_SET_DSCR_REGISTER(val)**, where value is a scalar of type integer*8
 - Set the DSCR value to the provided 64-bit value

The topic of DSCR is described from a processor perspective here:

- ▶ “Data prefetching using d-cache instructions and the Data Streams Control Register (DSCR)” on page 34.

7.3.5 Transactional memory (TM)

The XL and GCC compiler families provide common built-in functions for access to transactional memory (TM) features.

Information about the topic of TM, from the processor and OS perspectives, is available here:

- ▶ 2.2.4, “Transactional memory (TM)” on page 37 (*processor*)
- ▶ 4.2.4, “Transactional memory (TM)” on page 81 (*ALX*)
- ▶ 6.2.4, “Transactional memory (TM)” on page 113 (*Linux*)

7.4 IBM Feedback Directed Program Restructuring (FDPR)

Feedback Directed Program Restructuring (FDPR) is a feedback-based, directed, and post-link optimization tool.

7.4.1 Introduction

FDPR optimizes the executable binary file of a program by collecting information about the behavior of the program while the program is used for a typical workload, and then creates a new version of the program that is optimized for that workload. Both main executable and dynamically linked libraries (DLLs) are supported.

FDPR performs global optimizations at the level of the entire executable library, including statically linked library code. Because the executable library to be optimized by FDPR is not relinked, the compiler and linker conventions do not need to be preserved, thus allowing aggressive optimizations that are not available to optimizing compilers.

The main advantage that is provided by FDPR is the reduced footprint of both code and data, resulting in more effective cache usage. The principal optimizations of FDPR include global code reordering, global data reordering, function inlining, and loop unrolling, along with various tuning options tailored for the specific Power target. The effectiveness of the optimization depends largely on how representative the collected profile is regarding the true workload.

FDPR runs on both AIX and Linux and produces optimized code for all versions of the Power Architecture. POWER7 is its default target architecture.

Figure 7-1 shows how FDPR is used to optimize executable programs.

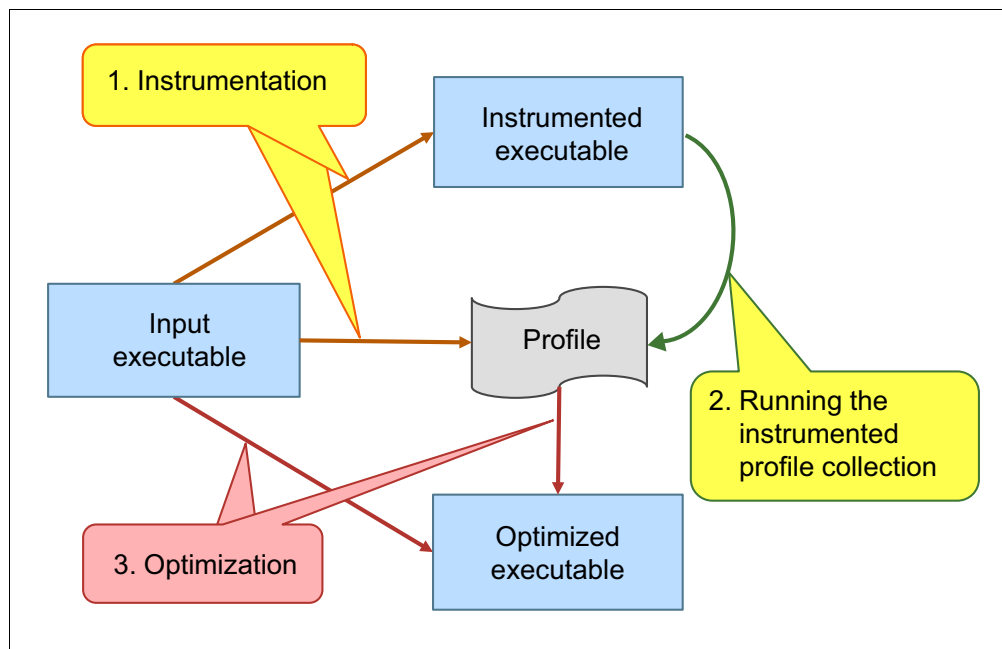


Figure 7-1 FDPR operation

FDPR builds an optimized executable program in three distinct phases:

1. Instrumentation (Yellow)
 - Creates an instrumented version of the input program and an empty profile file.
 - The input program can be an executable file or a dynamically linked shared library.
2. Profiling (Green)
 - Runs the instrumented program on a representative workload.
 - The profile file is filled with count data at run time.
3. Optimization (Red)

FDPR receives the original input program along with the filled profile file to create an optimized version of the input.

7.4.2 FDPR supported environments

FDPR is available on the following platforms:

- ▶ AIX and Power Systems: Part of the AIX 5L V5 operating system and higher for both 32-bit and 64-bit applications. For more information, see *AIX 5L Performance Tools Handbook*, SG24-6039:
- ▶ Software Development Toolkit for PowerLinux: Available for use through the IBM SDK for PowerLinux. Linux distributions of RHEL5 and above, and SLES10 and above are supported. For more information, see this website:

<http://www14.software.ibm.com/webapp/set2/sas/f/lopdiags/sdklop.html>

In these resources, detailed online help, including manuals, is provided for each of these environments.

7.4.3 Acceptable input formats

The input binary can be a main executable program or a shared library, originally written in any language (for example, C, C++, or Fortran), if it is statically compiled. Thus, Java byte code is not acceptable. Code that is written in assembly language is acceptable, but must follow the Power ABI convention. For more information, see *64-bit PowerPC ELF Application Binary Interface Supplement 1.9*, available here:

<http://refspecs.linuxfoundation.org/ELF/ppc64/PPC-elf64abi-1.9.pdf>

It is important that the file includes relocation information. Although this is the default in AIX, on Linux you must add **-Wl,-q**, or **-Wl,--emit-relocs** to the command used for linking the program (or **-q** if the **ld** command is used directly).

The input binary can include debug information. FDPR correctly processes line number information so that the optimized output can be debugged.

7.4.4 General operation

FDPR is started by running the **fdprpro** program as follows:

```
$ fdprpro -a action [-p] in -o out -f prof [opt ...]
```

The **action** indicates the specific processing that is requested. The most common ones are **instr** for the instrumentation step and **opt** for the optimization step.

The **in**, **out**, and **prof** indicate the input and output binary files and profile files.

FDPR comes also with a wrapper command, named **fdpr**, which performs the instrumentation, profiling, and optimization under one roof. Run **man fdpr** for more information about this wrapper.

Special input and output files

FDPR has a number of options that control input and output files. One option that controls the input files is **--ignored-function-list file (-ifl file)**.

In some cases, the structure of some functions confuses FDPR, which can result in bad code generation. The file that is specified by **--ignored-function-list file (-ifl file)** contains a list of functions that are considered unsafe for optimization. This configuration prevents the potential bad code generation that might otherwise occur.

In addition to the profile and the instrumented and output optimized files, FDPR can optionally produce various secondary files to help you understand the static and dynamic nature of the input binary program. These secondary files have the same base name as the output file and a special extension. The following options control important output files:

- ▶ **--disassemble_text (-d)** and **--dump-mapper (-dm)**: The **-d** option creates a disassembly of (the code segment) of the program (extension `.dis_text`). The disassembly is useful to understand the structure of program as analyzed or created by FDPR. The **-dm** option produces a mapping of basic-blocks from their original address to their address in the optimized code. This mapping can be used, for example, to understand how a specific piece of code was broken, or for user-specific post-processing tools.
- ▶ **--dump-ascii-profile (-dap)**: This option dumps the profile file in a human readable ASCII format (extension `.aprof`). The `.aprof` file is useful for manual inspection or user-defined post-processing of the collected profile.
- ▶ **--verbose *n* (-v *n*)**, **--print-inlined-funcs (-pif)**, and **--journal *file* (-j *file*)**: These options generate different analyses of the optimized file. **-v *n*** generates general and optimization-specific statistics (`.stat` extension). The amount of verbosity is set by *n*. Basic statistics are provided by **-v 1**. Optimization-specific statistics are added in level 2 and instruction mix in level 3. The list of inlining and inlined functions is produced with the **-pif** option (`.inl_list` extension). The **-j *file*** produces a journal of the main optimizations, in an XML format, with detailed information about each optimization site, including the corresponding source file and line information. This information can be used by GUI tools to display optimizations in the context of the source code.

Controlling output to the console

The amount of progress information that is printed to the console can be controlled by two options. The default progress information is as follows:

```
fdprpro (FDPR) Version vvv for Linux/POWER
fdprpro -a opt -O3 in -o out -f prof
> reading_exe ...
> adjusting_exe ...
> analyzing ...
> building_program_infrastructure ...
> building_profiling_cfg ...
> add_profiling ...
>> reading_profile ...
>> building_control_flow_transfer_profiling ...
> pre_reorder_optimizations ...
>> derat_optimization ...
...
```

This information might also be interspersed with warning and debugging messages. Use the **-quiet (-q)** option to avoid progress information. To limit the warning information, use the **-warning 1 (-w 1)** option.

7.4.5 Instrumentation and profiling

FDPR instrumentation is performed by running the following command:

```
$ fdprpro -a instr in [-o out] [-f prof] [opts...]
```

If **out** is not specified, the output file is `in.instr`. If the profile is not specified, `in.nprof` is used.

Two files are created, the instrumented program and an empty profile. The instrumented program (or shared library), when run on a representative workload, fills the profile with execution counts of nodes and edges of the binary control flow graph (CFG). A node in this CFG is a basic block (piece of code with single entry and exit points). An edge indicates a control transfer between two basic blocks through a branch (regular branch, call, or return instruction).

To run the instrumented program, use the same command parameters as with the original program. As indicated in 7.4.1, “Introduction” on page 149, the workload that is exercised during the instrumented run should be representative, making the optimization step more effective. Because of the instrumentation code, the program is slower.

Successive runs of the instrumented program accumulate the counts in the same profile. Similarly, if the instrumented program is a shared library, each time the shared library participates in a process, the corresponding profile is updated with added counts.

Profiling shared libraries

When the dynamic linker searches for and links a shared library during execution, it looks for the original name that is used to the command used for linking the program. To ensure that the instrumented library is run, ensure that the following items are true:

1. The instrumented library should have the same name as the original library. The user can rename the original or place the libraries in different folders.
2. The folder that contains the library must be in the library search path: `LIBPATH` on AIX and `LD_LIBRARY_PATH` on Linux.

Moving and renaming the profile file

The location of the profile file is specified in the instrumented program, as indicated by the `-f` option. However, the profile file might be moved, or if its original specification is relative, the real location can change before execution.

Use the `-fdir` option to set the profile directory if it is known at instrumentation time and is different from the one implied or specified by the `-f` option.

Use the `FDPR_PROF_DIR` environment variable to specify the profile directory if the profile file is not present in the relative or absolute location where it was created in the instrumentation step (or where specified originally by `-fdir`).

Use the `FDPR_PROF_NAME` environment variable to specify the profile file name if the profile file name changed.

Profile file descriptor (FD)

When the instrumented binary file is run, the profile file is mapped to shared memory. The process is using a default file descriptor (FD) number (1023 on Linux and 1999 on AIX) for the mapping. If the application uses this specific FD, an error can occur during the profiling phase because of this conflict of use. Use the `-fd` option to change the default FD used by FDPR:

```
$ fdprpro -a instr my_prog -fd <fd num>
```

The FD can also be controlled by using the `FDPR_PROF_FD` environment variable by changing the FD at run time:

```
$ export FDPR_PROF_FD=fd_num
```

FDPR can be used to profile several binary executable files in a single run of an application. If so, you must specify a different FD for each binary. Here is an example:

- ▶ `$ fdprpro -a instr in/libmy_lib1 -o out/libmy_lib1 -f out/libmy_lib1.prof -fd 1023`
- ▶ `$ fdprpro -a instr in/libmy_lib2 -o out/libmy_lib2 -f out/libmy_lib2.prof -fd 1022`

Because environment variables are global in nature, when profiling several binary files at the same time, use explicit instrumentation options (`-f`, `-fd`, and `-fdir`) to differentiate between the profiles rather than using the environment variables (`FDPR_PROF_FD` and `FDPR_PROF_NAME`).

Instrumentation stack

The instrumentation is using the stack for saving registers by dynamically allocating space on the stack at a default location below the current stack pointer. On AIX, this default is at offset -10240, and on Linux it is -1800. In some cases, especially in multi-threaded applications where the stack space is divided between the threads, following a deep calling sequence, the application can be quite close to the end of the stack, which can cause the application to fail. To allocate the instrumentation closer to the current stack pointer, use the `-iso` option:

```
$ fdprpro -a instr my_prog -iso -300
```

7.4.6 Optimization

The optimization step is performed by running the following command:

```
$ fdprpro -a opt in [-o out] -f prof [opts...]
```

If **out** is not specified, the output file is `in.fdpr`. No profile is provided by default. If **none** is specified or if the profile is empty, the resulting output binary file is not optimized.

Code reordering

Global code reordering works in two phases: making chains and reordering the chains.

The initial chains are sequentially ordered basic blocks, with branch conditions inverted where necessary, so that branches between the basic blocks are mostly not taken. This configuration makes instruction prefetching more efficient. Chains are terminated when the heat (that is, execution count) goes below a certain threshold relative to the initial heat.

The second phase orders chains by successively merging the more strongly linked two chains, based on how frequent the calls between the chains are. Combining chains crosses function boundaries. Thus, a function can be broken into multiple chunks in which different pieces of different functions are placed closely if there is a high frequency of call, branch, and return between them. This approach improves code locality and thus i-cache and page table efficiency.

You use the following options for code reordering:

- ▶ **--reorder-code (-RC)**: This component is the hard-working component of the global code reordering. Use **--rcaf** to determine the aggressiveness level:
 - 0: no change
 - 1: standard (default)
 - 2: most aggressive.

Use **--rcctf** to lower the threshold for terminating chains. Use **-pp** (preserve procedures) to preserve procedure boundaries in reordered code, and **-pc** to preserve CSECTS boundaries in reordered code. These two options limit global code reordering and might be requested for ease of debugging.

- **--branch-folding (-bf)** and **--branch-prediction (-bp)**: These options control important parts of the code reordering process. The **-bf** folds branch to branch into a single branch. The **-bp** sets the static branch prediction bit when taken or not taken statistics justify it.

Function inlining

FDPR performs function inlining of function bodies into their respective calling sites if the call site is selected by one of a number of user-selected filters:

- Dominant callers (**--selective-inlining (-si)**, **-sidf *f***, and **-siht *f***): The filter criteria here is that the site is dominant regarding other callers of the called function (the callee). It is controlled by two attributes. The **-sidf** option sets the domination percentage threshold (default 80). The **-siht** option further restricts the selection to functions hotter than the threshold, which is specified in percents relative to the average (default 100).
- Hot functions (**--inline-hot-functions *f* (-ihf *f*)**): This filter selects inlining for all call sites where the call is hotter than the heat threshold (in percent, relative to the average).
- Small functions (**--inline-small-functions *f* (-isf *f*)**): This filter selects for inlining all functions whose size, in bytes, is smaller than or equal to the parameter.
- Selective hot code (**--selective-hot-code-inline *f* (-shci *f*)**): The filter computes how much execution count is saved if the function is inlined at a call site and selects those sites where the relative saving is above the percentage.

De-virtualization

De-virtualization is addressed by the **--ptrgl-optimization (-pto)** option. It is full-blown call by a pointer mechanism (ptrgl) sets a new TOC anchor, loads the function address, moves it to the counter register (CTR), and jumps indirectly through the CTR. The **-pto** optimizes this mechanism in cases where there is few hot targets from a calling site. In terms of C++, it de-virtualizes the virtual method calls by calling the actual targets directly. The optimized code compares the address of the function descriptor, which is used for the indirect call, against the address of a hot candidate, as identified in the profile, and conditionally calls such target directly. If none of the hot targets match, the code invokes the original indirect call mechanism. The idea is that most of the time the conditional direct branches are run instead of the ptrgl mechanism. The impact of the optimization on performance depends heavily on the function call profile.

The following thresholds can help to tune the optimization and to adjust it to different workloads:

- Use **-ptoht *thres*** to set the frequency threshold for indirect calls that are to be optimized (*thres* can be 0 - 1, with 0.8 by default).
- Use **-ptos1 *n*** to set the limit of the number of hot functions to optimize in a given indirect call site (the default for *n* is 3).

Loop-unrolling

Most programs spend their time in loops. This statement is true regardless of the target architecture or application. FDPR has one option to control the unrolling optimization for loops: **--loop-unrolling *factor* (-lu *factor*)**.

FDPR optimizes loop using a technique called *loop-unrolling*. By unrolling a loop n times, the number of back branches is reduced n times, so code prefetch efficiency can be improved. The downside with loop-unrolling is code inflation, which results in increased code footprint and increased i-cache misses. Unlike traditional loop-unrolling, FDPR is able to mitigate this problem by unrolling only the hottest paths in the loop. The **factor** parameter determines the aggressiveness of the optimization. With **-03**, the optimization is invoked with **-1u 9**.

By default, loops are unrolled two times. Use **-1u factor** to change that default.

Architecture-specific optimizations

Here are some architecture-specific optimizations:

- ▶ **--machine tgt (-m tgt)**: FDPR optimizations include general optimizations that are based on a high-level program representation as a control and data flow, in addition to peephole optimizations, relying on different architecture features. Those optimizations can perform better when tuned for specific platforms. The **-m** flag allows the user to specify the target machine model when known in cases where the program is not intended for use on multiple target platforms. The default target is POWER7.
- ▶ **--align-code code (-A code)**: Optimizing the alignment and the placement of the code is crucial to the performance of the program. Correct alignment can improve instruction fetching and dispatching. The alignment algorithm in FDPR uses different techniques that are based on the target platform. Some techniques are generic for the Power Architecture, and others are considered dispatch rules of the specific machine model. If **code** is 1 (the default), FDPR applies a standard alignment algorithm that is adapted for the selected target machine (see **-m** in the previous bullet point). If **code** is 2, FDPR applies a more advanced version, using dispatch rules and other heuristics to decide how the program code chunks are placed relatively to i-cache sectors, again based on the selected target. A value of 0 disables the alignment algorithm.
- ▶ **--smt_mode (-smt)**: Grouping analysis depends on SMT mode. In general, FDPR tries to align code such that the number of runtime groups is minimized.

Function optimization

FDPR includes a number of function level optimizations that are based on detailed data flow analysis (DFA). With DFA, optimizations can determine the data that is contained in each register at each point in the function and whether this value is used later.

These are the function optimizations:

- ▶ **--killed-regs (-kr)**: A register is considered killed at a point (in the function) if its value is not used in any ensuing path. FDPR uses the Power ABI convention that defines which registers are non-volatile (NV) across function calls. NV registers that are used inside a function are saved in its prolog and restored in its epilog. The **-kr** optimization analyzes called functions that are looking for save and restore instructions of killed NV registers. If the register is killed at the calling site, then the save and restore instructions for this register are removed. The optimization considers all calls to this function, because an NV might be alive when the function is called. When needed, the optimization might also reassign (rename) registers at the calling side to ensure that an NV is indeed killed and can be optimized.
- ▶ **--hco-reschedule (-hr)**: The optimization analyzes the flow through hot basic blocks and looks for instructions that can be moved to dominating colder basic blocks (basic block b1 dominates b2 if all paths to b2 first go through b1). For example, an instruction that loads a constant to a register is a candidate for such motion.

- **--simplify-early-exit *factor* (-see *factor*)**: Sometimes a function starts with an early exit condition so that if the condition is met, the whole body of the function is ignored. If the condition is commonly taken, it makes sense to avoid saving the registers in the prolog and restoring them in the epilog. The **-see** optimization detects such a condition and provides a reduced epilog that restores only registers modified by computing the condition. If ***factor*** is 1, a more aggressive optimization is performed where the prolog is also optimized.

Peephole optimization

Peephole optimizations require a small context around the specific site in the code that is problematic. The more important ones that FDPR performs are **-rc1**, **-t1o**, **-las**, **-s1s**, and **-nop**:

- **--load-after-store (-las)**: In recent Power Architectures, when a load instruction from address A closely follows a store to that address, it can cause the load to be rejected. The instruction is then tried in a slower mode, which produces a large performance penalty. This behavior is also called *Load-Hit-Store (LHS)*. With the **-las** optimization, the load is pushed further from the store, thus avoiding the reject condition.
- **--toc-load-optimization (-t1o)**: The TOC is a data section in programs where pointers are kept to avoid the lengthy address computation at run time. Loading an address (a pointer) is a costly operation and FDPR is able to save on processing if the address is close enough to the TOC anchor (R2). In such cases, the load from TOC is replaced by an `addi Rt,R2,offset`, where `R2+offset = loaded address`. The optimization is performed after data is reordered so that commonly accessed data is placed closer to R2, increasing the potential of this optimization. A TOC is used in 32-bit and 64-bit programs on AIX, and in 64-bit programs on Power Systems running Linux. Linux 32-bit uses a GOT and this optimization is not relevant there.
- **--nop-removal (-nop)**: The compiler (or the linker) sometimes inserts no-operation (NOP) instructions in various places to create some necessary space in the instruction stream. The most common place is following a function call in code. Because the call might have modified the TOC anchor register (R2), the compiler inserts a load instruction that resets R2 to its correct value for the current function. Because FDPR has a global view of the program, the optimization can remove the NOP if the called function uses the same TOC (the TOC anchor is used in AIX and in Linux 64-bit).
- **--store-load-on-stack-opt (-s1s)**: POWER8 introduces new direct move instructions that allow for data movement between the general purpose registers and the floating point or floating vector registers (and additional new vector scalar instructions). This option looks for old code patterns that include such data movements using stores and loads to and from the stack, and replaces them with the new direct move instructions. The new instructions take less cycles and are expected to reduce pressure on other hardware resources, for example, store queue (STQ).
- **--xscpsgndp-to-xx1or (-xscpx)**: This optimization attempts to replace old **fmr** and **xscpsgndp** instructions, which take six cycles, with new **xx1or** efficient instructions, which take two cycles.
- **--remove-constant-load (-rc1)**: This optimization reduces the number of load instructions that are used to bring constant values into the registers. Many times, constants are loaded into a register using two load instructions, where the first instruction loads the address of the constant area from the table of contents (TOC), and the second instruction loads the constant value from the previously loaded address. FDPR tries to replace such patterns with faster computational instructions (no memory access).

Data reordering

The profile that is collected by FDPR provides important information about the running of branch instructions, thus enabling efficient code reordering. The profile does not provide this direct information whether to put specific objects one after the other. Nevertheless, FDPR is able to infer such placement by using the collected profile.

These are relevant options:

- ▶ **--reorder-data (-RD)**: This optimization reorders data by placing pointers and data closer to the TOC anchor, depending on their hotness. FDPR uses a heuristic where the hotness is computed as the total count of basic blocks where the pointer to the data was retrieved from the TOC.
- ▶ **--reduce-toc *thres* (-rt *thres*)**: The optimization removes from the TOC entries that are colder than the threshold. Their access, if any, is replaced by computing the address (see **-tlo** optimization in “Peephole optimization” on page 157). Typically, you use **-rt 0**, which removes only the entries that are never accessed.

Combination optimizations

FDPR has predefined optimization sets that provide a good starting point for performance tuning:

- ▶ **-0**: Performs code reordering (**-RC**) with branch prediction bit setting (**-bp**), branch folding (**-bf**), and NOOP instructions removal (**-nop**).
- ▶ **-02**: Adds to **-0** function de-virtualization (**-pto**), TOC-load optimization (**-tlo**), function inlining (**-isf 8**), and some function optimizations (**-hr**, **-see 0**, and **-kr**).
- ▶ **-03**: Switches on data reordering (**-RD** and **-rt 0**), loop-unrolling (**-lu**), more aggressive function optimization (**-see 1** and **-vro**), and employs more aggressive inlining (**-lro** and **-isf 12**). This set provides an aggressive but still stable set of optimizations that are beneficial for many benchmarks and applications.
- ▶ **-04**: Essentially turns on more aggressive inlining (**-sidf 50**, **-ihf 20**, and **-shci 90**). As a result, the number of branches is reduced, but at the cost of increasing code footprint. This option works well with large i-caches or with small to medium programs/threads.

7.5 Using the Advance Toolchain with IBM XLC and XLF

For XLC13 and XLF15, we have implemented a new feature into the existing `new_install` script, which is shipped with our Linux package.

Run this script with one option, and it will detect whether AT has been installed in the environment. If yes, it will automatically generate a config file with AT information specified, and meanwhile generate a new invoke named `xlc_at`, which uses the generated config file. Then you can use this `xlc_at` invoke to get the **XLC + AT** usage.

7.6 Related publications

The publications that are listed in this section are considered suitable for a more detailed discussion of the topics that are covered in this chapter:

- ▶ C/C++ Cafe Community:

<https://www.ibm.com/developerworks/community/groups/service/html/communityview?communityUuid=5894415f-be62-4bc0-81c5-3956e82276f3>

- ▶ *FDPR, Post-Link Optimization for Linux on Power*:

<https://www.ibm.com/developerworks/mydeveloperworks/groups/service/html/communityview?communityUuid=5a116d75-b560-4152-9113-7515fa73e67a>

- ▶ *Feedback Directed Program Restructuring (FDPR)*:

<https://www.research.ibm.com/haifa/projects/systems/cot/fdpr/>

- ▶ GCC online documentation:

- All versions: <http://gcc.gnu.org/onlinedocs/>
- Advance Toolchain 4.0: <http://gcc.gnu.org/onlinedocs/gcc-4.5.3/gcc/>
- Advance Toolchain 5.0: <http://gcc.gnu.org/onlinedocs/gcc-4.6.3/gcc/>
- Advance Toolchain 6.0: <http://gcc.gnu.org/onlinedocs/gcc-4.7.1/gcc/>
- Advance Toolchain 7.0: <http://gcc.gnu.org/onlinedocs/gcc-4.8.1/gcc/>

- ▶ XL compiler documentation:

- C and C++ compilers:

- *C and C++ Compilers family*:

<http://www.ibm.com/software/awdtools/xlcpp/>

- *Optimization and Programming Guide - XL C/C++ for AIX 13.1*:

<http://www.ibm.com/support/docview.wss?uid=swg27041856>

- Fortran compilers:

- *Fortran Compilers family*:

<http://www.ibm.com/software/awdtools/fortran/>

- *Optimization and Programming Guide - XL Fortran for AIX 15.1*:

<http://www.ibm.com/support/docview.wss?uid=swg27041932>



Java

This chapter describes the optimization and tuning of Java based applications that are running on a POWER8 processor-based server. It covers the following topics:

- ▶ 8.1, “Java levels” on page 162
- ▶ 8.2, “32-bit versus 64-bit Java” on page 162
- ▶ 8.3, “Memory and page size considerations” on page 163
- ▶ 8.4, “Java garbage collection tuning” on page 168
- ▶ 8.5, “Application scaling” on page 170
- ▶ 8.6, “Related publications” on page 175

8.1 Java levels

For POWER8, the preferred Java level is Java 7.1 or later. Java 7.1 is the only release of Java that is optimized for POWER8, and takes advantage of POWER8 specific hardware features for performance. Any version of Java 7.1 is acceptable, and later versions will contain additional POWER8 exploitation and tuning and are therefore preferred.

For POWER7, Java 7.1 is also preferred. However, it is acceptable to use Java 6 SR7 or later on POWER7. As of Java 6 SR7, the Java virtual machine (JVM) defaults to using 64 KB pages on AIX. Earlier versions defaulted to 4 KB pages, the default page size on AIX. (For more information, see “Tuning to capitalize on hardware performance features” on page 12, and 8.3.1, “Medium and large pages for Java heap and code cache” on page 167.)

The JIT compiler automatically detects on what platform it is running and generates binary code most suitable to, and performing best on, on that platform. Java 7.1 and later is able to recognize POWER8 and best use its hardware features.

8.2 32-bit versus 64-bit Java

64-bit applications that do not require large amounts of memory typically run slower than 32-bit applications. This situation occurs because of the larger data types, like 64-bit pointers instead of 32-bit pointers, which increase the demand on memory throughput.

The exception to this situation is when the processor architecture has more processor registers in 64-bit mode than in 32-bit mode and 32-bit application performance is negatively impacted by this configuration. Because of few registers, the demand on memory throughput can be higher in 32-bit mode than in 64-bit mode. In such a situation, running an application in 64-bit mode is required to achieve best performance.

The Power Architecture does not require running applications in 64-bit mode to achieve best performance because 32-bit and 64-bit modes have the same number of processor registers.

Consider the following items:

- ▶ Applications with a small memory requirement typically run faster as 32-bit applications than as 64-bit applications.
- ▶ 64-bit applications have a larger demand on memory because of the larger data types, such as pointers being 64-bit instead of 32-bit, which leads to the following circumstances:
 - The memory foot print increases because of the larger data types.
 - The memory alignment of application data contributes to memory demand.
 - More memory bandwidth is required.

For best performance, use 32-bit Java unless the memory requirement of the application requires running in 64-bit mode.

For more information about this topic, see 8.6, “Related publications” on page 175.

8.3 Memory and page size considerations

IBM Java can take advantage of medium (64 KB) and large (16 MB) page sizes that are supported by the current AIX versions and POWER processors. Using medium or large pages instead of the default 4 KB page size can improve application performance. The performance improvement of using medium or large pages is a result of a more efficient use of the hardware translation caches, which are used when you translate application page addresses to physical page addresses. Applications that are frequently accessing a vast amount of memory benefit most from using pages sizes that are larger than 4 KB.

Table 8-1 shows the hardware and software requirements for 4 KB, 64 KB, and 16 MB pages.

Table 8-1 Page sizes that are supported by AIX and Linux on Power Systems processors

Page size	Platform	Linux version	AIX version	Requires user configuration
4 KB	All	RHEL 5, SLES 10 and earlier	All	No
64 KB	POWER5+ or later	RHEL 6, SLES 11	AIX 5L V5.3 and later	No
16 MB	POWER4 or later	RHEL 5, SLES11	AIX 5L V5.3 and later	Yes

8.3.1 Medium and large pages for Java heap and code cache

Medium and large pages can be enabled for the Java heap and JIT code cache independently of other memory areas. IBM JVM supports at least three page sizes, depending on the platform:

- ▶ 4 KB (default)
- ▶ 64 KB
- ▶ 16 MB

Large pages, specifically 16 MB pages, do have some processing impact and are best suited for long running applications with large memory requirements. The **-Xlp64k** option provides many of the benefits of 16 MB pages with less impact and can be suitable for workloads that benefit from large pages but do not take full advantage of 16 MB pages.

Starting with IBM Java 6 SR7, the default page size is 64 KB.

Starting with IBM Java 7 SR4 (and Java 6.2.6 SR5), there are more command line options to specify pagesize for java heap and code cache. The **-Xlp:objectheap:pagesize=<size>** and **-Xlp:codecache:pagesize=<size>** options are supported. To obtain the large page sizes available and the current setting, use the **-verbose:sizes** option. Note that the current settings are the requested sizes and not the sizes obtained.

8.3.2 Configuring large pages for Java heap and code cache

In an AIX environment, to use large pages with Java requires both configuring the large pages and setting the **v_pinshm** tunable to a value of one by running **vmo**. The following example demonstrates how to dynamically configure 1 GB of 16 MB pages and set the **v_pinshm** tunable:

```
# vmo -o lpgg_regions=64 -o lpgg_size=16777216 -o v_pinshm=1
```

To permanently configure large pages, the **-r** option must be specified with the **vmo** command. Run **bosboot** to configure the large pages at boot time:

```
# vmo -r -o lpgg_regions=64 -o lpgg_size=16777216 -o v_pinshm=1
# bosboot -a
```

Non-root users must have the **CAP_BYPASS_RAC_VMM** capability on AIX enabled to use large pages. The system administrator can add this capability by running **chuser**:

```
# chuser capabilities=CAP_BYPASS_RAC_VMM,CAP_PROPAGATE <user_id>
```

On Linux, 1 GB of 16 MB pages are dynamically requested by running **sysctl**:

```
# sysctl -w vm.nr_hugepages=64
```

Use the following **cat** command to confirm the number of 16 MB pages that were actually allocated:

```
# cat /proc/sys/vm/nr_hugepages
```

Non-root users must be a member of a group on Linux that is enabled to use 16 MB pages. The system administrator can add this capability by running the following command:

```
# sysctl -w vm.hugetlb_shm_group=<group_id>
```

This allows all members of the specified **<group_id>** to use 16 MB pages.

To permanently configure the 16 MB pages on Linux, add the following lines to **/etc/sysctl.conf**:

```
vm.nr_hugepages = 64
vm.hugetlb_shm_group = <group_id>
```

Always check to ensure that Java is actually using large pages by testing using the **-verbosegc** option:

```
$ java -verbosegc -Xlp ...
```

In the **verbosegc** output, it will indicate both the requested page size and the actual page size that was used for the Java heap. Here is an example:

```
<attribute name="pageSize" value="0x1000000" />
<attribute name="requestedPageSize" value="0x1000000" />
```

This is necessary because Java will silently use a smaller page size if the 16 MB page environment is not correctly configured (that is, no warning or error message is issued).

8.3.3 Prefetching

Prefetching is an important strategy to reduce memory latency and take full advantage of on-chip caches. The **-XtlhPrefetch** option can be specified to enable aggressive prefetching of thread-local heap memory shortly before objects are allocated. This option ensures that the memory required for new objects that are allocated from the TLH is fetched into cache ahead of time if possible, reducing latency and increasing overall object allocation speed. This option can give noticeable gains on workloads that frequently allocate objects, such as transactional workloads.

8.3.4 Compressed references

For huge workloads, 64-bit IBM JVMs might be necessary to meet application needs. The 64-bit processes primarily offer a much larger address space, allowing for larger Java heaps, JIT code caches, and reducing the effects of memory fragmentation in the native heap. However, 64-bit processes also must deal with the increased processing impact. The impact comes from the increased memory usage and decreased cache usage. This impact is present with every object allocation, as each object must now be referred to with a 64-bit address rather than a 32-bit address.

To alleviate this impact, use the **-Xcompressedrefs** option. When this option is enabled, IBM JVM uses 32-bit references to objects instead of 64-bit references wherever possible. Object references are compressed and extracted as necessary at minimal cost. The need for compression and decompression is determined by the overall heap size and the platform that IBM JVM is running on; smaller heaps can do without compression and decompression, eliminating even this impact. To determine the compression and decompression impact for a heap size on a particular platform, run the following command:

```
java -Xcompressedrefs -verbose:gc -version ...
```

The resulting output has the following content:

```
<attribute name="compressedRefsDisplacement" value="0x0" />
<attribute name="compressedRefsShift" value="0x0" />
```

Values of 0 for the named attributes essentially indicate that no work must be done to convert between 32-bit and 64-bit references for the invocation. Under these circumstances, 64-bit IBM JVMs running with **-Xcompressedrefs** can reduce the impact of 64-bit addressing even more and achieve better performance.

With **-Xcompressedrefs**, the maximum size of the heap is much smaller than the theoretical maximum size allowed by a 64-bit IBM JVM, although greater than the maximum heap under a 32-bit IBM JVM. Currently, the maximum heap size with **-Xcompressedrefs** is around 31 GB on both AIX and Linux.

8.3.5 JIT code cache

JIT compilation is an important factor in optimizing performance. Because compilation is carried out at run time, it is complicated to estimate the size of the program or the number of compilations that are carried out. The JIT compiler has a cap on how much memory it can allocate at run time to store compiled code and for most of applications the default cap is more than sufficient.

However, certain programs, especially those programs that take advantage of certain language features, such as reflection, can produce a number of compilations and use up the allowed amount of code cache. After the limit of code cache is consumed, no more compilations are performed. This situation can have a negative impact on performance if the program begins to call many interpreted methods that cannot be compiled as a result. The `-Xjit:codetotal=<nnn>` (where *nnn* is a number in KB units) option can be used to specify the cap of the JIT code cache. The default is 64 MB or 128 MB for 32-bit and 64-bit IBM JVMs.

Another consideration is how the code caches are allocated. If they are allocated far apart from each other (more than 32 MB away), calls from one code cache to another carry higher processing impact. The `-Xcodecache<size>` option can be used to specify how large each allocation of code cache is. For example, `-Xcodecache4m` means 4 MB is allocated as code cache each time the JIT compiler needs a new one, until the cap is reached. Typically, there are multiple pieces (for example, 4) of code cache available at boot-up time to support multiple compilation threads. It is important to alter the default code cache size only if it is insufficient, as a large but empty code cache needlessly consumes resources.

`-Xcodecachetotal<size>` is the preferred option in Java 7.1. Java 7 SR6 and later, and Java 6 SR15 and later and are fully documented and supported.

Two techniques can be used to determine if the code cache allocation sizes or total limit must be altered. First, a Java core file can be produced by running `kill -3 <pid>` at the end/stable state of your application. The core file shows how many pieces of code cache are allocated. The active amount of code cache can be estimated by summing up all of the pieces.

For example, if 20 MB is needed to run the application, `-Xcodecache5m` (four pieces of 5 MB each) typically allocates 20 MB code caches at boot-up time, and they are likely close to each other and have better performance for cross-code cache calls. Second, to determine if the total code cache is sufficient, the `-Xjit:verbose` option can be used to print method names as they are compiled. If compilation fails because the limit of code cache is reached, an error to that effect is printed.

8.3.6 Shared classes

The IBM JVM supports class data sharing between multiple IBM JVMs. The `-Xshareclasses` option can be used to enable it, and the `-Xscmx<size>` option can be used to specify the maximum cache size of the stored data, where `<size>` can be `<nnn>K`, `<nnn>M`, or `<nnn>G` for sizes in KB, MB, or GB.

The shared class data is stored in a memory-mapped cache file on disk. Sharing reduces the overall virtual storage consumption when more than one IBM JVM shares a cache. Sharing also reduces the start time for a IBM JVM after the cache is created. The shared class cache is independent of any running IBM JVM and persists until it is deleted.

A shared cache can contain the following items:

- ▶ Bootstrap classes
- ▶ Application classes
- ▶ Metadata that describes the classes
- ▶ Ahead-of-time (AOT) compiled code

8.3.7 In-core Advanced Encryption Standard (AES) acceleration

Ensuring confidentiality through encryption is a computationally intensive aspect of workloads which is becoming increasingly important. POWER8 introduces in-core AES instructions that are compliant with the FIPS 197: AES Specification.

Starting with IBM Java 7.1, AES is accelerated using POWER8 in-core AES instructions by specifying `-Dcom.ibm.crypto.provider.doAESInHardware=true` on the JVM command line. In-core AES instructions can significantly increase speed, as compared with equivalent JIT-generated code.

8.3.8 Transactional memory (TM)

POWER8 Hardware Transaction Memory (HTM) is exploited by IBM JVM in two aspects, both of which are intended to be transparent to Java application programmers and JVM users:

- ▶ Transactional Lock Elision (TLE)
- ▶ Targeted class exploitation

Starting with IBM Java 7.1, and as long as HTM is enabled on the platform (AIX or Linux), this exploitation can occur transparently.

The JIT compiler automatically chooses particular Java synchronization blocks to transform into HTM regions. Only the blocks that are deemed to benefit from the transformation in terms of performance are chosen. When those blocks behave synergistically with HTM, application scalability and performance can be improved. Since the transformation is automatic, TLE is transparent to programmers and users. At the same time, certain classes, including `ConcurrentHashMap` and `ConcurrentLinkedQueue`, have been rewritten to take advantage of HTM in IBM JVM. These classes work on processors that do not support HTM, but they can take advantage of HTM running on POWER8 transparently.

However, application programmers can modify applications to take more advantage of TLE or HTM by targeting the applications specifically for POWER8. As explained in Chapter 2, `HashTable/HashMap/Map` generally behaves well with HTM and TLE. When the application is modified to use more of the classes mentioned (`ConcurrentHashMap` and `ConcurrentLinkedQueue`), the application is more likely to benefit from TLE.

Note: Using more of the classes mentioned (`ConcurrentHashMap` and `ConcurrentLinkedQueue`) may adversely affect performance on POWER7 or older processors that do not support HTM.

8.3.9 Runtime instrumentation

IBM Java 7 SR1 and later exploits the POWER8 event-based branching facility and enhanced performance monitoring unit (PMU) to enable runtime instrumentation of compiled code. Runtime instrumentation allows the JIT compiler to collect detailed performance information directly from the hardware, without any kernel or system call overhead, and in turn, use this information to further optimize compiled code. The PMU is the same unit that is currently used by external profiling tools, such as **tprof** and **hpmcount** on AIX, and **perf** and **OProfile** on Linux. POWER8 allows the PMU to be used for application self-profiling. JIT profiling and optimizations focus on collecting information that would otherwise be difficult to collect without hardware assistance and better utilization of cache and TLB resources, and on reducing function call overhead and branch mispredicts, among others.

8.4 Java garbage collection tuning

The IBM Java VM supports multiple garbage collection (GC) strategies to allow software developers an opportunity to prioritize various factors. Throughput, latency, and scaling are the main factors that are addressed by the different collection strategies. Understanding how an application behaves regarding allocation frequencies, required heap size, expected lifetime of objects, and other factors can make one or more of the non-default GC strategies preferable. The GC strategy can be specified with the **-Xgcpolicy:<policy>** option.

8.4.1 GC strategy: Optthruput

This strategy prioritizes throughput at the expense of maximum latency by waiting until the last possible time to do a GC. A global GC of the entire heap is performed, creating a longer pause time at the expense of latency. After GC is triggered, the GC stops all application threads and performs the three GC phases:

- ▶ Mark
- ▶ Sweep
- ▶ Compact (if necessary)

All phases are parallelized to perform GC as quickly as possible.

The optthruput strategy is the default in the original Java 6 that uses the V2.4 J9 IBM JVM.

8.4.2 GC strategy: Optavgpause

This strategy prioritizes latency and response time by performing the initial mark phase of GC concurrently with the execution of the application. The application is halted only for the sweep and compact phases, minimizing the total time that the application is paused. Performing the mark phase concurrently with the execution of the application might affect throughput, because the CPU time that would otherwise go to the application can be diverted to low priority GC threads to carry out the mark phase. This situation can be acceptable on machines with many processor cores and relatively few application threads, as idle processor cores can be put to good use otherwise.

8.4.3 GC strategy: Gencon

This strategy employs a generational GC scheme that attempts to deal with many varying workloads and memory usage patterns. In addition, gencon also uses concurrent marking to minimize pause times. The gencon strategy works by dividing the heap into two categories:

- ▶ New space
- ▶ Old space

The new space is dedicated to short-lived objects that are created frequently and unreferenced shortly thereafter. The old space is for long-lived objects that survived long enough to be promoted from the new space. This GC policy is suited to workloads that have many short-lived objects, such as transactional workloads, because GC in the new space (carried out by the *scavenger*) is cheaper per object overall than GC in the old space. By default, up to 25% of the heap is dedicated to the new space. The division between the new space and the old space can be controlled with the `-Xmn` option, which specifies the size of the new space; the remaining space is then designated as the old space. Alternatively, `-Xmns` and `-Xmnx` can be used to set the starting and maximum new space sizes if a non-constant new space size is wanted. For more information about constant versus non-constant heaps in general, see 8.4.5, “Optimal heap size” on page 170.

The gencon strategy is the default in the updated Java 6 that uses the V2.6 J9 VM, and in the later Java 7 and Java 7.1 versions.

8.4.4 GC strategy: Balanced

This strategy evens out pause times across GC operations that are based on the amount of work that is being generated. This strategy can be affected by object allocation rates, object survival rates, and fragmentation levels within the heap. This smoothing of pause times is a best effort rather than a real-time guarantee. A fundamental aspect of the balanced collector's architecture, which is critical to achieving its goals of reducing the impact of large collection times, is that it is a region-based garbage collector. A region is a clearly delineated portion of the Java object heap that categorizes how the associated memory is used and groups related objects together.

During the IBM JVM startup, the garbage collector divides the heap memory into equal-sized regions, and these region delineations remain static for the lifetime of the IBM JVM. Regions are the basic unit of GC and allocation operations. For example, when the heap is expanded or contracted, the memory that is committed or released corresponds to a number of regions.

Although the Java heap is a contiguous range of memory addresses, any region within that range can be committed or released as required. This situation enables the balanced collector to contract the heap more dynamically and aggressively than other garbage collectors, which typically require the committed portion of the heap to be contiguous. Java heap configuration for `-Xgcpolicy:balanced` strategy can be specified through the `-Xmn`, `-Xmx`, and `-Xms` options.

8.4.5 Optimal heap size

By default, the IBM JVM provides a considerably flexible heap configuration that allows the heap to grow and shrink dynamically in response to the needs of the application. This configuration allows the IBM JVM to claim only as much memory as necessary at any time, thus cooperating with other processes that are running on the system. The starting and maximum size of the heap can be specified with the `-Xms` and `-Xmx` options.

This flexibility comes at a cost, as the IBM JVM must request memory from the operating system whenever the heap must grow and return memory whenever it shrinks. This behavior can lead to various unwanted scenarios. If the application heap requirements oscillate, this situation can cause excessive heap growth and shrinkage.

If the IBM JVM is running on a dedicated machine, the processing impact of heap resizing can be eliminated by requesting a constant sized heap. This situation can be accomplished by setting `-Xms` equal to `-Xmx`. Choosing the correct size for the heap is highly important, as GC impact is directly proportional to the size of the heap. The heap must be large enough to satisfy the application's maximum memory requirements and contain extra space. The GC must work much harder when the heap is near full capacity because of fragmentation and other issues, so 20 - 30% of extra space above the maximum needs of the application can lower the overall GC impact.

If an application requires more flexibility than can be achieved with a constant sized heap, it might be beneficial to tune the sizing parameters for a dynamic heap. One of the most expensive GC events is *object allocation failure*. This failure occurs when there is not enough contiguous space in the current heap to satisfy the allocation, and results in a GC collection and a possible heap expansion. If the current heap size is less than the `-Xmx` size, the heap is expanded in response to the allocation failure if the amount of free space is below a certain threshold. Therefore, it is important to ensure that when an allocation fails, the heap is expanded to allow not only the failed allocation to succeed, but also many future allocations, or the next failed allocation might trigger yet another GC collection. This situation is known as *heap thrashing*.

The `-Xminf`, `-Xmaxf`, `-Xmine`, and `-Xmaxe` group of options can be used to affect when and how the GC resizes the heap. The `-Xminf<factor>` option (where factor is a real number 0 - 1) specifies the minimum free space in the heap; if the total free space falls below this factor, the heap is expanded. The `-Xmaxf<factor>` option specifies the maximum free space; if the total free space rises above this factor, the heap is shrunk. These options can be used to minimize heap thrashing and excessive resizing. The `-Xmine` and `-Xmaxe` options specify the minimum and maximum sizes to shrink and grow the heap by. These options can be used to ensure that the heap has enough free contiguous space to allow it to satisfy a reasonable number of allocations before failure.

Regardless of whether the heap size is constant, it should never be allowed to exceed the physical memory available to the process; otherwise, the operating system might have to swap data in and out of memory. An application's memory behavior can be determined by using various tools, including verbose GC logs. For more information about verbose GC logs and other tools, see "Java (either AIX or Linux)" on page 222.

8.5 Application scaling

Large workloads using many threads on multi-CPU machines face extra challenges regarding concurrency and scaling. In such cases, steps can be taken to decrease contention on shared resources and reduce the processing impact.

8.5.1 Choosing the correct SMT mode

AIX and Linux represent each SMT thread as a logical CPU. Therefore, the number of logical CPUs in an LPAR depends on the SMT mode. For example, an LPAR with four virtual processors that are running in SMT4 mode has 16 logical CPUs; an LPAR with that same number of virtual processors that are running in SMT2 mode has only eight logical CPUs.

Table 8-2 shows the number of SMT threads and logical CPUs available in ST, SMT2, and SMT4 modes.

Table 8-2 ST, SMT2, SMT4, and SMT8 modes - SMT threads and CPUs available

SMT mode	Number of SMT threads	Number of logical CPUs
ST	1	1
SMT2	2	2
SMT4	4	4
SMT8	8	8

The default SMT mode on POWER7 and later depends on the AIX version and the compatibility mode the processor cores are running with. Table 8-3 shows the default SMT modes.

Table 8-3 SMT mode on POWER8 is dependent upon AIX and compatibility mode

AIX version	Compatibility mode	Default SMT mode
AIX 7.1 TL3 SP3	POWER8	SMT4
AIX V6.1	POWER7	SMT4
AIX V6.1	POWER6/POWER6+	SMT2
AIX 5L V5.3	POWER6/POWER6+	SMT2

Most applications benefit from SMT. However, some applications do not scale with an increased number of logical CPUs on an SMT-enabled system. One way to address such an application scalability issue is to make a smaller LPAR, or use processor binding, as described in 8.5.2, “Using resource sets (RSETS)” on page 172.

Additionally, if you are in need of improved performance from your larger new system, there is a potential alternative. If your application semantics support it, you may be able to run multiple smaller instances of your application, each bound to exclusive processors/cores. In this way, the aggregate performance from the multiple instances of your application may be able to meet your performance expectations. This alternative is one of the WebSphere preferred practices. More information about selecting an appropriate SMT mode is available in , “Scalability challenges when moving from POWER5 or POWER6 to POWER7 or POWER8” on page 192.

For applications that might benefit from a lower SMT mode with fewer logical CPUs, experiment with using SMT2 or ST modes. For details from the processor and OS perspectives, see these sections:

- ▶ “SMT” on page 25 (*processor*)
- ▶ “Simultaneous Multithreading (SMT)” on page 65 (*AIX*)
- ▶ “SMT” on page 102 (*IBM i*)
- ▶ “Simultaneous multithreading (SMT)” on page 109 (*Linux*)

Java application scaling on Linux

Java applications scale better on Linux in some cases if the `sched_compat_yield` scheduler tunable is set to 1 by running the following command:

```
sysctl -w kernel.sched_compat_yield=1
```

Further information about this topic is available here:

- ▶ “Scalability considerations” on page 13 (*Linux*)
- ▶ “Simultaneous multithreading (SMT)” on page 109 (*Linux*)

8.5.2 Using resource sets (RSETS)

The use of RSETS in AIX and Linux environments follows.

AIX environment

In an AIX environment, resource sets (RSETS) allow specifying which logical CPUs an application can run on. They are useful when an application that does not scale beyond a certain number of logical CPUs is run on a large LPAR. For example, an application that scales up to eight logical CPUs but is run on an LPAR that has 64 logical CPUs.

See “The POWER8 processor and affinity performance effects” on page 14 for further information. An example is included in “Partition sizes and affinity” on page 15.

RSETS can be created with the `mkrset` command and attached to a process using the `attachrset` command. An alternative way is creating a resource set and attaching it to an application in a single step through the `execrset` command.

The following example demonstrates how to use `execrset` to create an RSET with CPUs 4 - 7 and run an application that is attached to it:

```
execrset -c 4-7 -e <application>
```

In addition to running the application attached to an RSET, set the `MEMORY_AFFINITY` environment variable to MCM to assure that the application’s private and shared memory is allocated from memory that is local to the logical CPUs of the RSET:

```
MEMORY_AFFINITY=MCM
```

In general, RSETs are created on core boundaries. For example, a partition with four POWER8 cores that are running in SMT4 mode has 16 logical CPUs. Create an RSET with four logical CPUs by selecting four SMT threads that belong to one core. Create an RSET with eight logical CPUs by selecting eight SMT threads that belong to two cores. The `smtctl` command can be used to determine which logical CPUs belong to which core, as shown in Example 8-1.

Example 8-1 Use the `smtctl` command to determine which logical CPUs belong to which core

```
# smtctl
This system is SMT capable.
This system supports up to 4 SMT threads per processor.
SMT is currently enabled.
SMT boot mode is not set.
SMT threads are bound to the same physical processor.
```

```

proc0 has 4 SMT threads.
Bind processor 0 is bound with proc0
Bind processor 1 is bound with proc0
Bind processor 2 is bound with proc0
Bind processor 3 is bound with proc0

proc4 has 4 SMT threads.
Bind processor 4 is bound with proc4
Bind processor 5 is bound with proc4
Bind processor 6 is bound with proc4
Bind processor 7 is bound with proc4

```

The **smtctl** output in Example 8-1 shows that the system is running in SMT4 mode with bind processors (logical CPU) 0 - 3 belonging to proc0 and bind processors 4 - 7 belonging to proc1. Create an RSET with four logical CPUs either for CPUs 0 - 3 or for CPUs 4 - 7.

To achieve the best performance with RSETs that are created across multiple cores, all cores of the RSET must be from the same chip and in the same scheduler resource allocation domain (SRAD). The **lssrad** command can be used to determine which logical CPUs belong to which SRAD, as shown in Example 8-2:

Example 8-2 Use the lssrad command to determine which logical CPUs belong to which SRAD

```

# lssrad -av
REF1 SRAD      MEM      CPU
0
      0 22397.25      0-31
1
      1 29801.75      32-63

```

The output in Example 8-2 shows a system that has two SRADs. CPUs 0 - 31 belong to the first SRAD, and CPUs 32 - 63 belong to the second SRAD. In this example, create an RSET with multiple cores either using the CPUs of the first or second SRAD.

Authority for RSETs: A user must have root authority or have **CAP_NUMA_ATTACH** capability to use RSETs.

Linux environment

In a Linux environment, the equivalent to **execrset** is the **taskset** command. The following example demonstrates how to use **taskset** to create a taskset with CPUs 4 - 7 and run an application that is attached to it:

```
Linux: taskset -c 4-7 <application>
```

There is no equivalent environment variable to **MEMORY_AFFINITY** on Linux; however, there is a command, **numactl**, that can accomplish the same task as does **MEMORY_AFFINITY** and also the **execrset** and **taskset** commands. Here is an example:

```
numactl [-l | --localalloc] -C 4-7 <application>
```

The **-l | --localalloc** option is analogous to **MEMORY_AFFINITY=MCM**.

8.5.3 Java lock reservation

Synchronization and locking are an important part of any multi-threaded application. Shared resources must be adequately protected by monitors to ensure correctness, even if some resources are only infrequently shared. If a resource is primarily accessed by a single thread at any time, that thread is frequently the only thread to acquire the monitor that is guarding the resource. In such cases, the cost of acquiring the monitor can be reduced by using the **-XlockReservation** option. With this option, it is assumed that the last thread to acquire the monitor is also likely to be the next thread to acquire it. The lock is, therefore, said to be reserved for that thread, minimizing its cost to acquire and release the monitor. This option is suited to workloads using many threads and many shared resources that are infrequently shared in practice.

8.5.4 Java GC threads

The GC used by the IBM JVM takes every opportunity to use parallelism on multi-CPU machines. All phases of the GC can be run in parallel with multiple helper threads dividing up the work to complete the task as quickly as possible. Depending on the GC strategy and heap size in use, it can be beneficial to adjust the number of threads that the GC uses. The number of GC threads can be specified with the **-Xgcthreads<number>** option. The default number of GC threads is generally equal to the number of logical processors on the partition, and it is usually not helpful to exceed this value. Reducing it, however, reduces the GC impact and might be wanted in some situations, such as when RSETs are used. The number of GC threads is capped at 64 starting in V2.6 J9 IBM JVM.

8.5.5 Java concurrent marking

The gencon policy combines concurrent marking with generational GC. If generational GC is wanted but the impact of concurrent marking, regarding both the impact of the marking thread and the extra book-keeping that is required when you allocate and manipulate objects, is not wanted, then concurrent marking can be disabled by using the **-Xconcurrentlevel0** option. This option is appropriate for workloads that benefit from the gencon policy for object allocation and lifetimes, but also require maximum throughput and minimal GC impact while the application threads are running.

In general, for both the gencon and optavgpause GC policies, concurrent marking can be tuned with the **-Xconcurrentlevel<number>** option, which specifies the ratio between the amounts of heap that is allocated and heap marked. The default value is 8. The number of low priority mark threads can be set with the **-Xconcurrentbackground<number>** option. By default, one thread is used for concurrent marking.

For more information about this topic, see 8.6, “Related publications” on page 175.

8.6 Related publications

The publications that are listed in this section are considered suitable for a more detailed discussion of the topics that are covered in this chapter:

- *Java Performance on POWER7 – Best Practice*, found here:

<http://public.dhe.ibm.com/common/ssi/ecm/en/pow03066usen/POW03066USEN.PDF>

- *Java Performance on POWER7*, found here:

<https://www.ibm.com/developerworks/wikis/display/LinuxP/Java+Performance+on+POWER7>

- *Top 10 64-bit IBM WebSphere Application Server FAQ*, found here:

ftp://public.dhe.ibm.com/software/webservers/appserv/WAS_64-bit_FAQ.pdf



DB2

This chapter describes the optimization and tuning of DB2 running on POWER processor-based servers. It covers the following topics:

- ▶ 9.1, “DB2 and the POWER processor” on page 178
- ▶ 9.2, “Taking advantage of the POWER processor” on page 178
- ▶ 9.3, “Capitalizing on the compilers and optimization tools for POWER” on page 181
- ▶ 9.4, “Capitalizing on POWER virtualization” on page 182
- ▶ 9.5, “Capitalizing on the AIX system libraries” on page 183
- ▶ 9.6, “Capitalizing on performance tooling” on page 185
- ▶ 9.7, “Conclusion” on page 186
- ▶ 9.8, “Related publications” on page 186

9.1 DB2 and the POWER processor

IBM DB2 is well-positioned to take full advantage of the POWER architecture. In this chapter, we are referring to DB2 Version 10, including DB2 10.1, DB2 10.5 and all subsequent updates. References to *POWER* are to the POWER7, POWER7+ and the new POWER8 processors. DB2 offers many capabilities that are tailored to use POWER features. The DB2 self-tuning memory manager (STMM) feature is one of many features that can help DB2 workloads efficiently consolidate on POWER Systems.

Additionally, DB2 is one of the most optimized software applications on POWER. During the DB2 development cycles, we always evaluate new POWER processor capabilities and ensure that we test and characterize the performance of DB2 on the latest POWER servers. So far, in the earlier chapters in this book, you read detailed descriptions about most of these POWER guidelines and technologies. The focus of this chapter is to showcase how IBM DB2 uses various POWER features and preferred practices from this guide during its own software development cycle, which is done to maximize performance on the Power Architecture. General DB2 tuning and preferred practices of DB2 are covered extensively in many other places, some of which are listed at in 9.8, “Related publications” on page 186.

Most of the POWER exploitation capabilities of DB2 extend to the POWER7, POWER7+ and POWER8 processors. A number of the new POWER8 capabilities are evolutionary in nature, and no new externals are required in DB2, just verification and some adjustment of the internal data structures and algorithms in DB2 to take full advantage of the new capabilities. Similarly, DB2 evolves, so maximizing on the benefits of POWER in DB2 10.1 has been extended and enhanced in DB2 10.5.

However, there are also new capabilities in DB2 10.5 that take advantage of never before used POWER processor capabilities. For example, the new columnar-in-memory analytic processing capability known as BLU Acceleration makes use of the POWER VSX engine on all of the POWER7 and later processors.

9.2 Taking advantage of the POWER processor

Methods for taking advantage of the inherent power of the POWER processor include affinization, page size, decimal arithmetics, and the usage of SMT priorities for internal lock implementation. New in DB2 10.5 is the ability to use Single Instruction Multiple Data (SIMD) processing with the VSX engine.

9.2.1 Affinization

New to DB2 10.1 is a much easier way to achieve affinization on POWER7 and later systems through the DB2 registry variable `DB2_RESOURCE_POLICY`. In general, this variable defines which operating system resources are available for DB2 databases or assigns specific resources to a particular database object. A typical example is to define a resource policy that restricts a DB2 database to run only on a specific set of processor cores.

On POWER7 and later systems running AIX V6.1 Technology Level (TL) 5 or higher, this variable can be set to AUTOMATIC. This feature is a new one introduced in DB2 10.1 and can result in enhanced query performance on some workloads. With this setting, the DB2 database system automatically detects the Power hardware topology and computes the best way to assign engine dispatchable units (EDUs) to various hardware modules. The goal is to determine the most efficient way to share memory between multiple EDUs that need access to the same regions of memory. This setting is intended for larger POWER7 and later systems with 16 or more cores. It is best to run a performance analysis of the workload before and after you set this variable to AUTOMATIC to validate the performance improvement.

The following example demonstrates how to set the registry variable to AUTOMATIC using the **db2set** command and then starting the DB2 database manager:

- ▶ **db2set DB2_RESOURCE_POLICY=AUTOMATIC**
- ▶ **db2start**

In DB2 10.1, DB2_RESOURCE_POLICY uses Scheduler Resource Allocation Domain Identifier (SRADID) instead of resource set (RSET) attachments to identify resources for the database.

For more information about other usages of DB2_RESOURCE_POLICY other memory-related DB2 registry variables, see Chapter 2, “AIX configuration”, in *Best Practices for DB2 on AIX 6.1 for POWER Systems*, SG24-7821.

9.2.2 Page sizes

Physical objects on disks, such as tables and indexes, are stored in pages. DB2 supports 4 KB, 8 KB, 16 KB, and 32 KB page sizes. During the processing of such objects, they are brought into DB2 buffer pools in main memory. The default AIX page size is 4 KB, but other page sizes are available. To achieve increased performance on Power Systems, DB2 10.1 by default uses 64 KB, which is a medium page size.

Large page size

For some workloads, particularly ones that require intensive memory access, there are performance benefits in using large page size support on AIX. However, certain drawbacks must be considered. When large page size support is enabled through DB2, all the memory that is set for large pages is pinned. It is possible that too much memory is allocated for large pages and not enough for 4 KB pages, which can result in heavy paging activities. Furthermore, enabling large pages prevents the STMM from automatically tuning overall database memory consumption. Consider using this variable only for well-defined workloads that have a relatively static database memory requirement.

POWER7 and later large page size support can be enabled by setting the DB2 registry variable DB2_LARGE_PAGE_MEM.

Here are the steps to enable large page support in DB2 database system on AIX operating systems:¹

1. Configure AIX server for large pages support by running **vmo**:

```
vmo -r -o lpgg_size=LargePageSize -o lpgg_regions=LargePages
```

LargePageSize is the size in bytes of the hardware-supported large pages, and LargePages specifies the number of large pages to reserve.
2. Run **bosboot** to pick up the changes (made by running **vmo**) for the next system boot.
3. After reboot, run **vmo** to enable memory pinning:

```
vmo -o v_pinshm=1
```
4. Set the DB2_LARGE_PAGE_MEM registry variable by running **db2set**, then start the DB2 database manager by running **db2start**:
 - **db2set DB2_LARGE_PAGE_MEM=DB**
 - **db2start**

9.2.3 Decimal arithmetics

DB2 uses the hardware DFP unit in IBM POWER6 and later processors in its implementation of decimal-encoded formats and arithmetics. One example of a data type that uses this hardware support is the DECFLOAT data type that is introduced in DB2 9.5. This decimal-floating point data type supports business applications that require exact decimal values, with precision of 16 or 34 digits. When the DECFLOAT data type is used for a DB2 database that is on a POWER6 or later processor, the native hardware support for decimal arithmetic is used. In comparison to other platforms, where such business operations can be achieved only through software emulation, applications that run on POWER6 or later can use the hardware support to gain performance improvements.

DECFLOAT: The Data Type of the Future describes this topic in more detail. The paper is available here:

<http://www.ibm.com/developerworks/data/library/techarticle/dm-0801chainani/>

9.2.4 Using SMT priorities for internal lock implementation

DB2 uses SMT and hardware priorities in its internal lock implementation. Internal locks are short duration locks that are required to ensure consistency of various values in highly concurrent applications such as DB2. In certain cases, it is beneficial to prioritize different DB2 agent threads to maximize system resource utilization.

For more information about this topic, see 9.8, “Related publications” on page 186.

9.2.5 SIMD

DB2 10.5 with BLU Acceleration uses SIMD processing to further speed up analytic query processing. On POWER processors, this is called the VSX engine. The VSX engine has been part of the POWER architecture for a time, was improved in POWER7, and has received further improvements in POWER8.

¹ *Enabling large page support (AIX)* (for DB2 Version 10.1 for Linux, UNIX, and Windows), available here:
<http://pic.dhe.ibm.com/infocenter/db2luw/v10r1/index.jsp?topic=%2Fcom.ibm.db2.luw.admin.dboobj.doc%2Fdoc%2Ft0010405.html>

SIMD instructions are low-level CPU instructions that enable you to perform the same operation on multiple data points at the same time.

DB2 10.5 with BLU Acceleration auto-detects whether it is running on an SIMD-enabled CPU, and automatically uses SIMD to effectively multiply the power of the CPU. In particular, BLU Acceleration can use a single SIMD instruction to get results from multiple data elements.

Figure 9-1 is an example of a scan operation that involves a predicate evaluation.

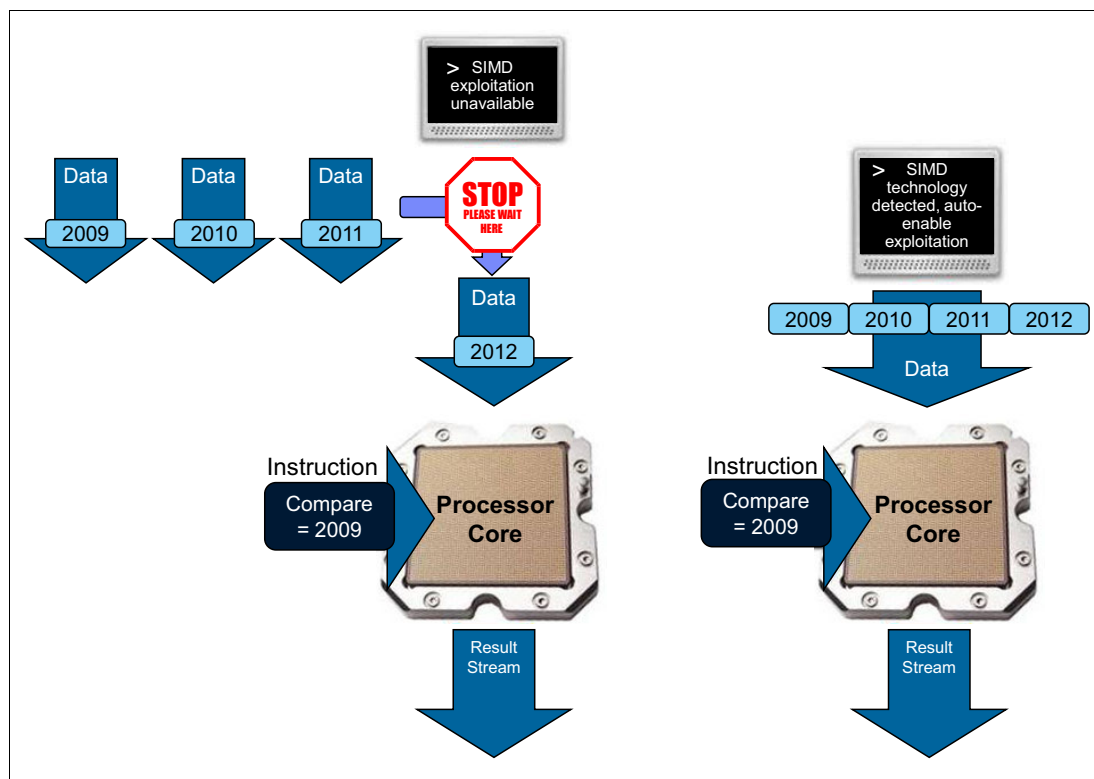


Figure 9-1 Comparing predicate evaluation with and without SIMD on POWER using DB2 10.5 with BLU Acceleration

The left side of Figure 9-1 illustrates a typical operation, namely, that each data element (or column value) is evaluated, one after another. The right side of the figure shows how BLU Acceleration processes four columns at a time by the use of SIMD. Think of it as CPU power, multiplied by four. Although Figure 9-1 shows a predicate evaluation, BLU Acceleration can also take advantage of SIMD processing for join operations, arithmetic, and more.

9.3 Capitalizing on the compilers and optimization tools for POWER

DB2 is built by using the latest IBM XL C/C++ compiler available at the start of our development cycle. For DB2 10.1, we use the Version 11 compiler, and for DB2 10.5, we use the Version 12 compiler. We use various compiler optimization flags along with optimization techniques based on the common three steps of software profiling:

- ▶ Application preparation/instrumentation
- ▶ Application profiling
- ▶ Application optimization

9.3.1 Whole-program analysis and profile-based optimizations

On the AIX platform, whole-program analysis (IPA) and profile-based optimization (PDF) compiler options are used to optimize DB2 using a set of customer representative workloads. This technique produces a highly optimized DB2 executable file that is targeted at making the best usage of the Power Architecture.

For further information, see “Whole-program analysis (IPA)” on page 139 and “Optimization that is based on Profile Directed Feedback” on page 139.

9.3.2 Feedback directed program restructuring (FDPR)

In addition to IPA and PDF optimizations using IBM XL C/C++ compiler, a post-link optimization step provides further performance improvement on the AIX platform. The particular tool that is used is IBM Feedback Directed Program Restructuring (FDPR). Similar to IPA and PDF, a set of DB2 customer representative workloads is employed in this step, and IBM FDPR-Pro profiles and ultimately creates an optimized version of the DB2 product. For further information, see 7.4, “IBM Feedback Directed Program Restructuring (FDPR)” on page 149.

For more information about this topic, see 9.8, “Related publications” on page 186.

9.4 Capitalizing on POWER virtualization

DB2 supports and fully draws upon the virtualization technologies that are provided by the POWER7 and later architectures. These technologies include PowerVM for AIX and Linux and System Workload Partitioning (WPAR) for AIX. Many of the DB2 performance preferred practices for non-virtualized environment also extend to a virtualized environment.

Furthermore, DB2 offers IBM SubCapacity Licensing, which enables customers to more effectively consolidate their infrastructure and reduce their overall total cost of ownership (TCO). DB2 also provides a flexible software licensing model that supports advanced virtualization capabilities, such as shared processor pools, Micro-Partitioning, virtual machines, and dynamic reallocation of resources. To support this type of licensing model, a tool is provided that allows customers to track and manage their own software license usage.

9.4.1 DB2 virtualization

DB2 is engineered to take advantage of the many benefits of virtualization on POWER and therefore allows various types of workload to be deployed in a virtualized environment. One key DB2 feature that enables workloads to run efficiently in virtualized environments is STMM. STMM is designed to automatically adjust the values of several memory configuration parameters in DB2. When enabled, it dynamically evaluates and redistributes available memory resources among the buffer pools, lock memory, package cache, and sort memory to maximize performance. The changes are applied dynamically and can simplify the task of manual configuration of memory parameters. This feature is useful in a virtualized environment because STMM can respond to dynamic changes in partition memory allocation.

By default, most DB2 parameters are set to automatic to enable STMM. As a preferred practice, leave the **instance_memory** parameter and other memory parameters as automatic, especially when you are running in a virtualized environment because DB2 is designed to allow STMM to look for available memory in the system when **instance_memory** is set to automatic.

DB2 also supports the PowerVM Live Partition Mobility (LPM) feature when virtual I/O is configured. LPM allows an active database to be moved from a system with limited memory to one with more memory without disrupting the operating system or applications. When coupling dynamic LPAR (DLPAR) with STMM, the newly migrated database can automatically adjust to the additional memory resource for better performance.

DB2 Virtualization, SG24-7805, describes in considerable detail the concept of DB2 virtualization, in addition to setup, configuration, and management of DB2 on IBM Power Systems with PowerVM technology. That book follows many of the preferred practices for Power virtualization and has a list of preferred practices for DB2 on PowerVM.

9.4.2 DB2 in an AIX workload partition

DB2 supports product installation on system WPARs. DB2 can be installed either within a local file system on a system WPAR or in a global environment under either the `/usr` or `/opt` directory with each instance created on the local WPARs. In both cases, each DB2 instance is only visible and managed by the system WPAR that it is created in. If DB2 is installed in a global environment, different instances on different WPARs share the globally installed DB2 copy to improve i-cache efficiency and memory usage. WPAR mobility is also supported where a DB2 instance that is running on a system WPAR can migrate to a remote WPAR on a different physical machine.

There are certain restrictions and considerations to keep in mind when you install DB2 in a global environment:

- ▶ Certain DB2 installation features cannot be installed on a system WPAR. These features are IBM Tivoli® System Automation for Multiplatforms and IBM Data Studio Administration Console.
- ▶ When you uninstall a DB2 copy in a global environment, all associated instances must be dropped or updated to another DB2 copy and its corresponding system WPARs must be active.
- ▶ When you apply fix packs to a DB2 copy in a global environment, all associated instances must be stopped and its corresponding system WPARs must be active.

For information about installing a DB2 copy on a WPAR, see Chapter 8, “Workload Partitioning”, in *Best Practices for DB2 on AIX 6.1 for POWER Systems*, SG24-7821.

For more information about this topic, see 9.8, “Related publications” on page 186.

9.5 Capitalizing on the AIX system libraries

This section describes methods for capitalizing on the AIX system libraries.

9.5.1 Using the `thread_post_many` API

DB2 uses `thread_wait` and `thread_post_many` to improve the efficiency of DB2 threads running on multi-processor Power Systems. DB2 takes advantage of the `thread_post_many` function. The availability of such an API on AIX directly impacts the efficiency of DB2 processing, as it allows for waking many EDUs with a single function call, which in other operating systems requires many individual function calls (typically as many as the number of EDUs being woken up).

9.5.2 File systems

DB2 uses most of the advanced features within the AIX file systems. These features include Direct I/O (DIO), Concurrent I/O (CIO), Asynchronous I/O, and I/O Completion Ports (IOCP).

Non-buffered I/O

By default, DB2 uses CIO or DIO for newly created table space containers because non-buffered I/O provides more efficient underlying storage access over buffered I/O on most workloads, with most of the benefit that is realized by bypassing the file system cache. Non-buffered I/O is configured through the `NO FILE SYSTEM CACHING` clause of the table space definition. To maximize the benefits of non-buffered I/O, a correct buffer pool size is essential. This size can be achieved by using STMM to tune the buffer pool sizes. (The default buffer pool is always tuned by STMM, but user created buffer pools must specify the **automatic** keyword for the size to allow STMM to tune them.) When STMM is enabled, it automatically adjusts the buffer pool size for optimal performance.

For file systems that support CIO, such as AIX JFS2, DB2 automatically uses this I/O method because of its performance benefits over DIO.

The DB2 log file by default uses DIO, which brings similar performance benefits as avoiding file system cache for table spaces.

Asynchronous I/O

In general, DB2 users cannot explicitly choose synchronous or asynchronous I/O. However, to improve the overall response time of the database system, minimizing synchronous I/O is preferred and can be achieved through correct database tuning. Consider the following items:

- ▶ Synchronous read I/O can occur when a DB2 agent needs a page that is not in the buffer pool to process an SQL statement. In addition, a synchronous write I/O can occur if no clean pages are available in the buffer pool to make room to bring another page from disk into that buffer pool. This situation can be minimized by having sufficiently large buffer pools or setting the buffer pool size to automatic to allow STMM to find its optimal size, in addition to tuning the page cleaning (by using the `chnpggs_thresh` database parameter).
- ▶ Not all pages read into buffer pools are done synchronously. Depending on the SQL statement, DB2 can prefetch pages of data into buffer pools through asynchronous I/O. When prefetching is enabled, two parallel activities occur during query processing: data processing and data page I/O. The latter is done through the I/O servers that wait for prefetch requests from the former. These prefetch requests contain a description of the I/O that must satisfy the query. The number of I/O servers for a database is specified through the `num_ioservers` configuration parameter. By default, this parameter is automatically tuned during database startup.

For more information about how to monitor and tune AIO for DB2, see *Best Practices for DB2 on AIX 6.1 for POWER Systems*, SG24-7821.

I/O Completion Port (IOCP)

Configure the AIX I/O completion port for performance purposes, even though it is not mandatory, as part of the DB2 version 10 installation process. For more information, see *Configuring IOCP (AIX)*, available here:

<http://pic.dhe.ibm.com/infocenter/db2luw/v10r1/index.jsp?topic=/com.ibm.db2.luw.admin.perf.doc/doc/t0054518.html>

After IOCP is configured on AIX, then DB2, by default, capitalizes on this feature for all asynchronous I/O requests. With IOCP configured, AIO server processes from the AIX operating system manage the I/O requests by processing many requests in the most optimal way for the system.

For more information about this topic, see 9.8, “Related publications” on page 186.

9.6 Capitalizing on performance tooling

Correct performance tooling is crucial for maximizing DB2 performance. Zoning in on potential performance bottlenecks is impossible without a strong performance tool set, such as the ones on Power Systems.

9.6.1 High-level investigation

During the general analysis of any performance investigation, the identification of the system resource bottlenecks is key to determining the root cause of the bottleneck. System resource bottlenecks can be classified into several categories, such as CPU bound, IO bound, network bound, or excessive idling, all of which can be identified with AIX system commands.

9.6.2 Low-level investigation

Various system level tools are essential in drilling down to find a potential root cause for the type of the bottlenecks that are listed in 9.6, “Capitalizing on performance tooling” on page 185. Profiling tools are especially invaluable for identifying CPU bound issues and are available on AIX and Linux platforms for POWER.

AIX tprof

tprof is a powerful profiling tool on AIX and Linux platforms that does program counter-sampling in clock interrupts. It can work on any binary without recompilation and is a great tool for codepath analysis.

For instructions about using *The tprof command*, visit this website:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.prftools/doc/prftools/tprofcommand.htm>

AIX tprof microprofiling

Beyond the high-level **tprof** profiling, DB2 also uses the microprofiling option of **tprof** during development. Microprofiling allows DB2 to perform instruction level profiling to attributes the CPU time spent on source program lines.

Linux OProfile

OProfile is a system profiling tool, similar in nature to **tprof**, that is popular on the Linux platform. **OProfile** uses hardware counters to provide functional level profiling in both the kernel and user space. Like **tprof**, this tool is useful during DB2 development for codepath analysis.

For more information about this topic, see 9.8, “Related publications” on page 186.

9.7 Conclusion

DB2 is positioned to capitalize on many Power features to maximize the ROI of the full IBM stack. During the entire DB2 development cycle, there is a targeted effort to take advantage of Power features and ensure that the highest level of optimization is employed on this platform. With every new Power generation, DB2 ensures that the key features are supported and brought into play at Power launch time, by working on such features well in advance of general availability. This type of targeted effort ensures that DB2 is at the forefront of optimization for Power applications.

9.8 Related publications

The publications that are listed in this section are considered suitable for a more detailed discussion of the topics that are covered in this chapter:

- ▶ *Best Practices for DB2 on AIX 6.1 for Power Systems*, SG24-7821
- ▶ *Best practices for DB2 for Linux, UNIX, and Windows*, found here:
<http://www.ibm.com/developerworks/data/bestpractices/db2luw/>
- ▶ *DB2 database products in a workload partition (AIX)*, found here:
<http://pic.dhe.ibm.com/infocenter/db2luw/v10r1/index.jsp?topic=/com.ibm.db2.luw.qb.server.doc/doc/c0053344.html>
- ▶ DB2 performance registry variables, including **DB2_LOGGER_NON_BUFFERED_IO** and **DB2_USE_IOCP**, are described in *Performance variables*, found here:
<http://pic.dhe.ibm.com/infocenter/db2luw/v10r1/index.jsp?topic=/com.ibm.db2.luw.admin.regvars.doc/doc/r0005665.html>
- ▶ *DB2 Version 10.1 for Linux, UNIX, and Windows, Performance variables* describes DB2 performance registry variables, including **DB2_RESOURCE_POLICY** and **DB2_LARGE_PAGE_MEM**:
<http://pic.dhe.ibm.com/infocenter/db2luw/v10r1/index.jsp?topic=/com.ibm.db2.luw.admin.regvars.doc/doc/r0005665.html>
- ▶ *DB2 Virtualization*, SG24-7805
- ▶ *DECFLOAT: The data type of the future*, found here:
<http://www.ibm.com/developerworks/data/library/techarticle/dm-0801chainani/>
- ▶ *DECFLOAT scalar function*, (for DB2 Version 10.1 for Linux, UNIX, and Windows) found here:
<http://pic.dhe.ibm.com/infocenter/db2luw/v10r1/index.jsp?topic=/com.ibm.db2.luw.sql.ref.doc/doc/r0050508.html>
- ▶ *Feedback Directed Program Restructuring (FDPR)*, found here:
<https://www.research.ibm.com/haifa/projects/systems/cot/fdpr/>

- ▶ *FDPR-Pro - Usage: Feedback Directed Program Restructuring*, found here:
http://www.research.ibm.com/haifa/projects/systems/cot/fdpr/papers/fdpr_pro_usage_cs.pdf
- ▶ IBM DB2 Version 10.1 Information Center, found here:
<http://pic.dhe.ibm.com/infocenter/db2luw/v10r1/index.jsp?topic=/com.ibm.db2.luw.welcome.doc/doc/welcome.html>
- ▶ *Smashing performance with OProfile*, found here:
<http://www.ibm.com/developerworks/library/l-oprof/>
- ▶ *tprof Command*, found here:
<http://pic.dhe.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.cmds/doc/aixcmds5/tprof.htm>



WebSphere Application Server

This chapter describes the optimization and tuning of the POWER8 processor-based server running WebSphere Application Server. It covers the following topics:

- ▶ 10.1, “IBM WebSphere on Power Systems” on page 190
- ▶ 10.2, “Performance and functional considerations” on page 190

10.1 IBM WebSphere on Power Systems

This chapter is intended to provide you with performance and functional considerations for running WebSphere Application Server middleware on Power Systems. It primarily describes POWER7 and POWER8. Even though WebSphere Application Server is designed to run on many operating systems and platforms, some specific capabilities of Power Systems are used by WebSphere Application Server as a part of platform optimization efforts.

10.2 Performance and functional considerations

The intent of this chapter is to explain WebSphere Application Server installation, deployment, and migration topics when WebSphere Application Server is running on Power Systems. This chapter describes preferred practices for performance when you run enterprise Java applications on Power Systems. It also highlights some of the known WebSphere Application Server topics and solutions for Power Systems.

10.2.1 Installation

As there are multiple versions of WebSphere Application Server, there are also multiple versions of AIX supported by POWER7 and POWER8. Table 10-1 shows some of the installation considerations. We suggest using the most currently available code, including the latest installation binaries. For the most current AIX installation and configuration details, see 4.4.1, “AIX preferred practices that are applicable to all Power Systems generations” on page 96.

Important: If running on POWER7, use the versions of WebSphere Application Server that run in *POWER7 mode with performance enhancements*. Similarly, for POWER8, use the versions of WebSphere Application Server that run in *POWER8 mode with performance enhancements*.

Table 10-1 Installation considerations

Consideration	Associated website	Information about website
IBM WebSphere Application Server support on POWER7 hardware	http://www.ibm.com/support/docview.wss?uid=swg21422150	Various fix pack levels and 64-bit considerations for running in POWER7 mode

10.2.2 Deployment

When you start the WebSphere Application Server, there is an option to bind the Java processors to specific CPU processor cores to circumvent the operating system scheduler to send the work to available processors in the pool. In certain cases, using RSETs and binding the JVM to stay within core/socket boundaries improves the performance. Table 10-2 lists some of the deployment considerations.

Table 10-2 Deployment considerations

Consideration	Associated website	Information about website
<i>Workload partitioning (WPAR) in AIX V6.1</i>	http://www.ibm.com/developerworks/aix/library/au-wpar61aix/	Determining when it is useful to move from LPAR deployment to WPAR deployment
<i>Troubleshooting and performance analysis of different applications in versioned WPARs</i>	http://www.ibm.com/developerworks/aix/library/au-wpars/	The benefits of moving from old hardware to the new POWER7 hardware in the form of versioned WPARs

Processor affinity benefits for WebSphere applications

When an application is running on top of WebSphere Application Server is deployed on a large LPAR, it might not use all the cores in that LPAR, resulting in less than optimum application performance. If this situation occurs, performance improvements to these applications can be obtained by binding the application server to certain cores. This task can be accomplished by creating the resource sets and attaching them to the application server that is running **excerset**. For an example of using the **taskset** and **numactl** commands in a Linux environment, see “Partition sizes and affinity” on page 15.

10.2.3 Performance

When you run WebSphere Application Server on POWER7 and POWER8 systems, end-to-end performance depends on many subsystems. This included the network, memory, disk, and CPU subsystems of POWER7 and POWER8; a crucial consideration is Java configuration and tuning. Topology also plays a major role in the performance of the enterprise application that is being deployed. The architecture of the application must be considered when you determine the best deployment topology. Table 10-3 includes links to preferred practices documents, which target each of these major areas.

Table 10-3 Performance considerations

Consideration	Associated website	Information provided
<i>Java Performance on POWER7 - Best practice</i>	http://public.dhe.ibm.com/common/ssi/ecm/en/pow03066usen/POW03066USEN.PDF	This white paper highlights key preferred practices for all Java applications that are running on Power Systems and simultaneous multithreading (SMT) considerations when you are migrating from POWER5 or POWER6 to POWER7.
<i>Optimizing AIX 7 network performance: Part 1, Network overview - Monitoring the hardware</i>	http://www.ibm.com/developerworks/aix/library/au-aix7networkoptimize1/index.html	This three-part white paper reviews AIX V7.1 networking and includes suggestions for achieving the best network performance.
<i>Optimizing AIX V7 memory performance: Part 1, Memory overview and tuning memory parameters</i>	http://www.ibm.com/developerworks/aix/library/au-aix7memoryoptimize1/index.html	Memory optimization is essential for running WebSphere Application Server faster on POWER7.
<i>Optimizing AIX V7 performance: Part 2, Monitoring logical volumes and analyzing the results</i>	http://www.ibm.com/developerworks/aix/library/au-aix7optimize2/index.html	Optimizing the disk and troubleshooting the I/O bottlenecks is crucial for I/O-intensive applications.

WebSphere channel framework degradation on POWER7

Certain applications that run on WebSphere Application Server on POWER7 can experience performance degradation because of asynchronous I/O (AIO). AIO can be disabled to improve the performance of these applications. For instructions about how to accomplish this task, see *Disabling AIO (Asynchronous Input/Output) native transport in WebSphere Application Server*, available here:

<http://www.ibm.com/support/docview.wss?uid=swg21366862>

Scalability challenges when moving from POWER5 or POWER6 to POWER7 or POWER8

By default, POWER7 runs in SMT4 mode, and POWER8 runs in SMT4 (AIX) or SMT8 (Linux, IBMi) modes. As such, there are either 4 or 8 hardware threads (logical CPUs) per core that provide tremendous concurrency for applications. If the enterprise applications are migrated to POWER7 or POWER8 from an earlier version of POWER hardware (POWER5 or POWER6), you might experience scalability issues, because the default SMT mode on POWER8 is SMT4 (AIX) or SMT8 (Linux, IBMi), and POWER7 is SMT4, but on POWER5 and POWER6, the default is SMT and SMT2 mode, respectively. As some of these applications might not be designed for the massive parallelism of POWER7 or POWER8, performance and scalability can be improved by using smaller partitions or processor binding. Processor binding is described in “Processor affinity benefits for WebSphere applications” on page 191.

Memory affinity benefits for WebSphere applications

In addition to the processor affinity described in “Processor affinity benefits for WebSphere applications”, applications can benefit from avoiding remote memory accesses by setting the environment variable `MEMORY_AFFINITY` to `MCM`. This variable allocates application private and shared memory from processor local memory.

These three tuning techniques (SMT scalability, CPU affinity, and memory affinity) can improve the performance of WebSphere Application Server on POWER7 and POWER8 systems. For an example of using the `taskset` and `numactl` commands in a Linux environment, see “Partition sizes and affinity” on page 15. See also “Processor affinity benefits for WebSphere applications” on page 191.

More information about these topics is provided in *Java Performance on POWER7 - Best practice*, found here:

<http://public.dhe.ibm.com/common/ssi/ecm/en/pow03066usen/POW03066USEN.PDF>

10.2.4 Performance analysis, problem determination, and diagnostic tests

Resources for addressing issues regarding performance analysis, problem determination, and diagnostic tests are listed in Table 10-4.

Table 10-4 Performance analysis and problem determination

Consideration	Associated website	Information provided
<i>Java Performance Advisor (JPA)</i>	https://www.ibm.com/developerworks/wikis/display/WikiPtype/Java+Performance+Advisor	The JPA tool provides suggestions for improving the performance of Java/WebSphere Application Server applications that are running on Power Systems.
<i>The performance detective: Where does it hurt?</i>	http://www.ibm.com/developerworks/aix/library/au-performancedetective/index.html	Describes how to isolate performance problems.
<i>MustGather: Performance, hang, or high CPU issues with WebSphere Application Server on AIX</i>	http://www.ibm.com/support/docview.wss?uid=swg21052641	MustGather assists with collecting the data necessary to diagnose and resolve issues with hanging or CPU usage issues.

Further details about addressing performance analysis are included in these references:

- ▶ 8.2, “32-bit versus 64-bit Java” on page 162
- ▶ 8.3, “Memory and page size considerations” on page 163
- ▶ 8.4, “Java garbage collection tuning” on page 168
- ▶ 8.5, “Application scaling” on page 170



A

Analyzing malloc usage under AIX

This appendix describes the optimization and tuning of the memory usage of an application by using the AIX malloc subroutine. It covers the following topics:

- ▶ “Introduction” on page 196
- ▶ “How to collect malloc usage information” on page 197

Introduction

There is a simple methodology on AIX to collect useful information about how an application uses the C heap. That information can then be used to choose and tune the appropriate malloc settings. The type of information that typically must be collected is:

- ▶ The distribution of malloc allocation sizes that are used by an application, which shows if AIX **MALLOCOPTIONS**, such as pool and buckets, are expected to perform well. This information can be used to fine-tune bucket sizes.
- ▶ The steady state size of the heap, which shows how to size the pool option.

Additional information about thread counts, malloc usage per thread, and so on, can be useful, but the information that is presented here presents a basic view.

This discussion does not apply to the watson2 allocator (see “Memory allocators” on page 86), which autonomically adjusts to the memory usage of an application and does not require specific tuning.

How to collect malloc usage information

To discover the distribution of allocation sizes, set the following environment variable:

```
export MALLOCOPTIONS=buckets,bucket_statistics:stdout
```

Run an application. When the application completes, a summary of the malloc activity is output. Example A-1 shows a sample output from a simple test program.

Example: A-1 Output from a simple test program

```
=====
Malloc buckets statistical summary
=====
Configuration values:
  Number of buckets: 16
  Bucket sizing factor: 32
  Blocks per bucket: 1024
Allocation request totals:
  Buckets allocator:    118870654
  Default allocator:    343383
  Total for process:    119214037
Allocation requests by bucket
Bucket      Maximum      Number of
Number      Block Size   Allocations
-----
   0           32      104906782
   1           64       9658271
   2           96      1838903
   3          128       880723
   4          160       300990
   5          192       422310
   6          224       143923
   7          256       126939
   8          288       157459
   9          320        72162
  10          352        87108
  11          384        56136
  12          416        63137
  13          448        66160
  14          480        45571
  15          512        44080
Allocation requests by heap
Heap      Buckets      Default
Number    Allocator    Allocator
-----
   0      118870654      343383
```

This environment variable causes the program to produce a histogram of allocation sizes when it terminates. The number of allocation requests satisfied by the default allocator indicates the fraction of requests that are too large for the buckets allocator (larger than 512 bytes, in this example). By modifying some of the malloc buckets configuration options, you can, for example, obtain more details about larger allocation sizes.

To discover the steady state size of the heap, set the following environment variable:

```
export MALLOCDDEBUG=log
```

Run an application to a steady state point, attach it by running **dbx**, and then run **malloc**. Example A-2 shows a sample output.

Example: A-2 Sample output from the malloc subroutine

```
(dbx) malloc
The following options are enabled:
    Implementation Algorithm..... Default Allocator (Yorktown)
    Malloc Log
        Stack Depth..... 4
Statistical Report on the Malloc Subsystem:
    Heap 0
        heap lock held by..... pthread ID 0x20023358
        bytes acquired from sbrk()..... 5309664
        bytes in the freespace tree..... 334032
        bytes held by the user..... 4975632
        allocations currently active..... 76102
        allocations since process start.. 20999785
The Process Heap
    Initial process brk value..... 0x20013850
    current process brk value..... 0x214924c0
    sbrk()s called by malloc..... 78
```

The bytes held by the user value indicates how much heap space is allocated. By stopping multiple times when you run **dbx** and then running **malloc**, you can get a good estimate of the heap space needed by the application.

For more information, see *System memory allocation using the malloc subsystem*, available here:

http://publib.boulder.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.genprogc/doc/genprogc/sys_mem_alloc.htm



Performance tooling and empirical performance analysis

This appendix describes the optimization and tuning of the POWER8 processor-based server from the perspective of performance tooling and empirical performance analysis. It covers the following topics:

- ▶ “Introduction” on page 200
- ▶ “Performance advisors” on page 200
- ▶ “Power Virtualization Performance (PowerVP)” on page 206
- ▶ “AIX” on page 206
- ▶ “Linux” on page 215
- ▶ “Java (either AIX or Linux)” on page 222

Introduction

This appendix includes a general description about performance advisors, and descriptions specific to the three performance advisors that are referenced in this book:

- ▶ AIX
- ▶ Linux
- ▶ Java (either AIX or Linux)

Performance advisors

IBM developed four new performance advisors that empower users to address their own performance issues to best use their Power Systems server. These performance advisors can be run by a broad class of users.

The first three of these advisors are tools that run and analyze the configuration of a system and the software that is running on it. They also provide advice about the performance implications of the current configuration and suggestions for improvement. These three advisors are documented in “Expert system advisors” on page 200.

The fourth advisor is part of the IBM Rational Developer for Power Systems Software™. It is a component of an integrated development environment (IDE), which provides a set of features for performance tuning of C and C++ applications on AIX and Linux. That advisor is documented in “Rational Performance Advisor” on page 205.

Expert system advisors

The expert system advisors are three new tools that are developed by IBM. What is unique about these applications is that they collect and interpret performance data. In one step, they collect performance metrics, analyze data, and provide a one-page visual report. This report summarizes the performance health of the environment, and includes instructions for alleviating detected problems. The performance advisors produce advice that is based on the expertise of IBM performance analysts, and IBM documented preferred practices. These expert systems focus on AIX Partition Virtualization, VIOS, and Java performance.

All of the advisors follow the same reporting format, which is a single page XML file you can use to quickly assess conditions by visually inspecting the report and looking at the descriptive icons, as shown in Figure B-1.






ICON	DEFINITION
	Informative: Context relevant data helpful in making adjustments.
	Optimal: Current condition likely to deliver best performance.
	Warning: Current condition deviates from best practices. Opportunity likely exists for better performance.
	Critical: Current condition likely causing negative impacts.
	Investigate: Further investigation or information required by user to determine if observation is impacting performance.

Figure B-1 Descriptive icons in expert system advisors (AIX Partition Virtualization, VIOS Advisor, and Java Performance Advisor)

The XML reports generated by all of the advisors are interactive. If a problem is detected, three pieces of information are shared with the user:

1. What is this?
This section explains why a particular topic was monitored, and provides a definition of the performance metric or setting.
2. Why is it Important?
This report entry explains why the topic is relevant and how it impacts performance.
3. How do I modify?
Instructions for addressing the problem are listed in this section.

VIOS Performance Advisor

The VIOS Performance Advisor provides guidance about various aspects of VIOS, including:

- ▶ CPU
- ▶ Shared processing pool
- ▶ Memory
- ▶ Fibre Channel performance
- ▶ Disk I/O subsystem
- ▶ Shared Ethernet adapter

The output is presented on a single page, and copies of the report can be saved, making it easy to document the settings and performance of VIOS over time. The goal of the advisor is for you to be able to self-assess the health of your VIOS and act to attain optimal performance.

Figure B-2 is a panel of the VIOS Performance Advisor, focusing on the FC adapter section of the report, which attempts to guide the user in determining if any of the FC ports are being saturated, and, if so, to what extent. An investigate image was displayed next to the idle FC port to confirm that the idle adapter port is intentional and because of an administrative configuration design choice.

The *VIOS Advisor* can be found here:

<https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/Power%20Systems/page/VIOS%20Advisor>

Figure B-2 shows a panel from the VIOS Advisor.







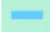


VIOS - DISK ADAPTERS							
	Name	Measured Value	Recommended Value	First Observed	Last Observed	Risk 1=lowest 5=highest	Impact 1=lowest 5=highest
	FC Adapter Count	2	-	06/12 20:32:25	-	n/a	n/a
	FC Avg IOps	avg: 7134 iops @ 4KB	 show/hide details	06/12 20:32:25	06/12 20:37:25	n/a	n/a
	FC Avg IOps(fcs2)	idle	-	06/12 20:32:25	06/12 20:37:25	n/a	n/a
	FC Avg IOps(fcs3)	avg: 7134 iops @ 4KB peak: 7251 iops @ 4KB	-	06/12 20:32:25	06/12 20:37:25	n/a	n/a
	FC Adapter Utilization		 show/hide details	-	-	n/a	n/a
	FC Adapter Utilization(fcs2)	idle	-	06/12 20:32:44	06/12 20:37:16	4	4
	FC Adapter Utilization(fcs3)	high:14.3% util. avg:14.1%	-	06/12 20:32:44	06/12 20:37:16	4	4

Figure B-2 The VIOS Advisor

Virtualization Performance Advisor

The Virtualization Performance Advisor provides guidance for various aspects of a logical partition (LPAR), both dedicated and shared, including these:

- ▶ LPAR physical memory domain allocation
- ▶ Physical CPU entitlement and virtual CPU optimization
- ▶ SMT effectiveness
- ▶ Processor folding effectiveness
- ▶ Shared processing pool
- ▶ Memory optimization
- ▶ Physical Fibre Channel adapter optimization
- ▶ Virtual disk I/O optimization (Virtual small computer system interface (vSCSI) and N_Port ID Virtualization (NPIV))

The output is presented on a single page, and copies of the report can be saved, making it easy for the user to document the settings and performance of their LPAR over time. The goal of the advisor is for the user to be able to self-assess the health of their LPAR and act to attain optimal performance.

Figure B-3 is a snapshot of the LPAR performance advisor, focusing on the LPAR optimization section of the report, which applies virtualization preferred practice guidance to LPAR configuration, resource usage of the LPAR, and shared processor pool, and determines if the LPAR configuration is optimized. If the advisor finds that the LPAR configuration is not optimal for the workload, it guides the user in determining the best possible configuration. The *LPAR Performance Advisor* can be found here:

https://www.ibm.com/developerworks/community/blogs/simplyaix/entry/lpar_performance_advisor?lang=en

LPAR PROCESSOR OPTIMIZATION										
	Name	Current Value				Recommended Value	First Observed	Last Observed	Risk 1=lowest 5=highest	Impact 1=lowest 5=highest
✖	Lpar Placement Optimization	Placement				Memory is allocated from multiple domains, containing memory in one or fewer domains will improve performance. However, changing this allocation requires assistance from IBM. Also, refer to P7 Virtualization Performance best practice document.	Tue Mar 6 20:14:28 2012	-	NA	NA
		Global Domain	Chip Domain	Memory	CPU					
		0	0	63228.560	0-11					
		0	6	60756.000						
		1	1	63246.000	12-15					
		1	2	63246.000	16-19					
		2	3	63223.000	20-23					
		2	4	63478.690	24-27					
		3	5	61503.000	28-31					
3	7	52539.000								
Only memory assigned from 0 Global Domain 6 Chip Domain Only memory assigned from 3 Global Domain 7 Chip Domain										
✖		Local Memory access - 0.00 and Remote memory access - 1210.84 Remote memory access is high. Local Memory access - 0.00 and Distant memory access - 9275.13 Distant memory access is high.				Memory is allocated from multiple domains, containing memory in one or fewer domains will improve performance. However, changing this allocation requires assistance from IBM. Also, refer to P7 Virtualization Performance best practice document.	Tue Mar 6 20:14:28 2012	-	NA	NA
✔	SMT Effectiveness	SMT-4 is set				-	Tue Mar 6 20:14:28 2012	-	NA	NA
✔	Virtual Processor Folding Optimization	Virtual Processor Folding Threshold - 49%				Rerun when lpar is busy	Tue Mar 6 20:14:33 2012	-	NA	NA

Figure B-3 LPAR Virtualization Advisor

Java Performance Advisor

The Java Performance Advisor provides advice to improve performance of a stand-alone Java or WebSphere Application Server application that is running on an AIX machine. The guidance that is provided is categorized into four groups as follows:

- ▶ Hardware and LPAR-related parameters: Processor sharing, SMT levels, memory, and so on
- ▶ AIX specific tunables: Process RSET, TCP buffers, memory affinity, and so on
- ▶ JVM tunables: Heap sizing, garbage collection (GC) policy, page size, and so on
- ▶ WebSphere Application Server related settings for a WebSphere Application Server process

The guidance is based on Java tuning preferred practices. The criteria that are used to determine the guidance include the relative importance of the Java application, machine usage (test and production), and the user's expertise level.

Figure B-4 is a snapshot of Java and WebSphere Application Server suggestions from a sample run, indicating the best JVM optimization and WebSphere Application Server settings for better results, as per Java preferred practices. Details about the metrics can be obtained by expanding each of the metrics. The output of the run is a simple XML file that can be viewed by using the supplied XSL viewer and any browser. The *Java Performance Advisor (JPA)* can be found here:

<https://www.ibm.com/developerworks/wikis/display/WikiPtype/Java+Performance+Advisor>













Java					
	Name	Current Value	Recommended Value	Risk 1=lowest 5=highest	Impact 1=lowest 5=highest
	JVM Version	1.6.0 SR2	More Details...	4	4
	JVM Type	64 bit	32 bit	4	4
	Initial Heap Size	100 MB	400 MB to 1.5625 GB	2	3
	Maximum Heap Size	1.5625 GB	1.5625 GB	5	5
	JVM Debug	Off	Off	1	5
	Verbose Class Loading	Off	Off	1	2
	Verbose Garbage Collection	Off	On	1	1
WebSphere					
	Name	Current Value	Recommended Value	Risk 1=lowest 5=highest	Impact 1=lowest 5=highest
	WebSphere Version	7.0.0.0	More Details...	3	4
	WebSphere PMI	On	Off	3	5
	Session Time Out	30 Minutes	5 Minutes to 30 Minutes	3	1
	Minimum Web Container Threads	50 Threads	10 Threads to 72 Threads	3	3
	Maximum Web Container Threads	50 Threads	24 Threads to 144 Threads	4	5

Figure B-4 Java Performance Advisor

Rational Performance Advisor

IBM Rational Developer for Power Systems Software IDE V8.5 introduces a new component that is called Performance Advisor, which provides a rich set of features for performance tuning C and C++ applications on IBM AIX and IBM PowerLinux systems. Although not directly related to the tooling described in “Expert system advisors” on page 200, Rational Performance Advisor has the same goal of helping users to best use Power hardware with tooling that offers simple collection, management, and analysis of performance data.

Performance Advisor gathers data from several sources. The raw application performance data comes from the same expert-level **tprof** and **OProfile** CPU profilers described in “AIX” on page 206 and “Linux” on page 215, and other low-level operating system tools. The debug information that is generated by the compiler allows this data to be matched back to the original source code. XLC compilers can generate XML report files that provide information about optimizations that were performed during compilation. Finally, the application build and runtime systems are analyzed to determine whether there are any potential environmental problems.

All of this data is automatically gathered, correlated, analyzed, and presented in a way that is quick to access and easy to understand (Figure B-5).

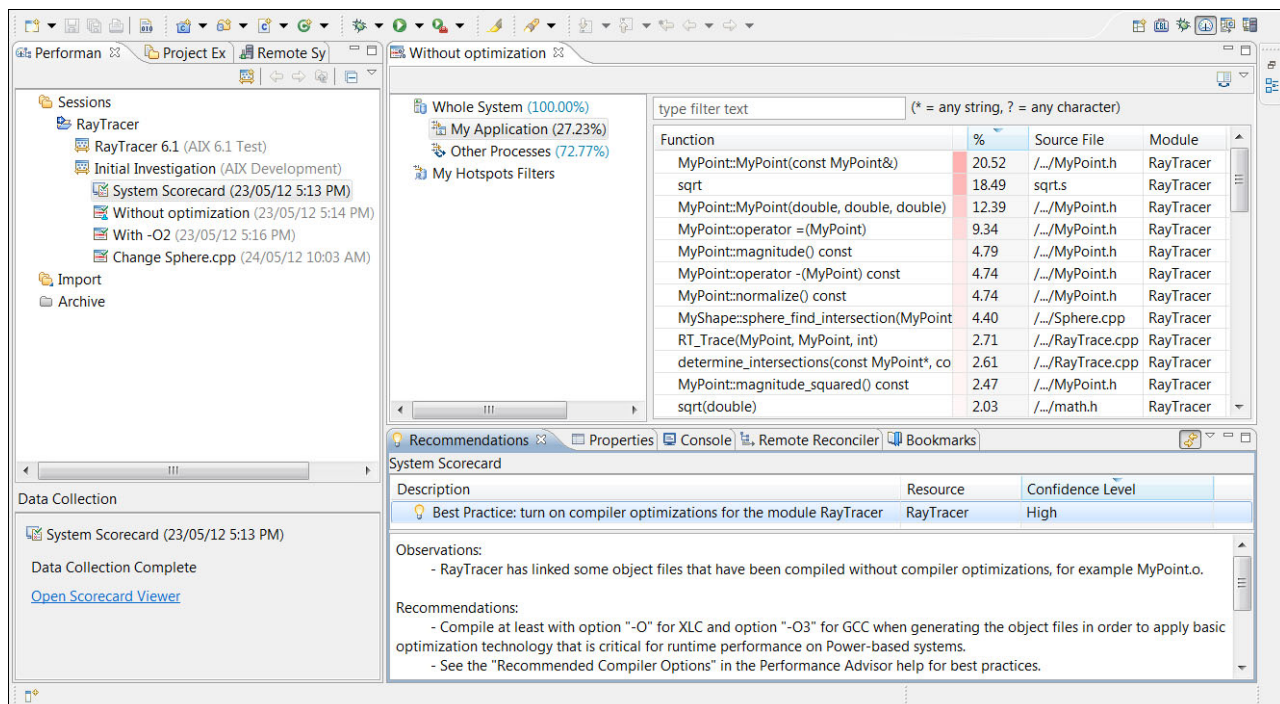


Figure B-5 Rational Performance Advisor

Key features include these:

- ▶ Performance Explorer organizes your performance tuning sessions and data.
- ▶ System Scorecard reports on your Power build and runtime environments.
- ▶ Hotspots Browser shows CPU profiling results for your application and its functions.
- ▶ Hotspots Comparison Browser compares runs for regression analysis or fix verification.
- ▶ The Performance Source Viewer and Outline view gives precise line-level profiling results.
- ▶ Invocations Browser displays dynamic call information from your application
- ▶ The Recommendations view offers expert-system guidance.

More information about Rational Performance Advisor, including a trial download, can be found in *Rational Developer family*, available here:

<http://www.ibm.com/software/rational/products/rdp/>

Power Virtualization Performance (PowerVP)

IBM PowerVP is a performance monitoring solution that provides detailed and real-time information about virtualized workloads that are running on Power Systems. PowerVP is a licensed program that is offered as part of PowerVM Enterprise Edition, but is also available separately for clients without PowerVM Enterprise Edition. You can use PowerVP to understand how virtual workloads use resources, to analyze performance bottlenecks, and to make informed choices about resource allocation and virtualized machine placement. PowerVP version 1.1.2 supports the POWER8 hardware.

The PowerVP tool:

- ▶ Monitors the performance of an entire system (or frame).
- ▶ Is supported on AIX, IBM i, Linux, and Virtual I/O Server operating systems.
- ▶ Provides a GUI for monitoring virtualized workloads.
- ▶ Includes a system-level monitoring agent that collects data from the PowerVM hypervisor, which provides a complete view of virtualized machines that are running on the server.
- ▶ Displays the data that is collected at the system level, at the hardware node level, and at the partition level. You can optimize performance by using the PowerVP performance metrics, which provide information about balancing and improving affinity and application efficiency.
- ▶ Provides an illustration of the Power Systems hardware topology along with resource usage metrics.
- ▶ Provides a mapping between real and virtual processor resources.
- ▶ Provides a recording feature for storing performance information with digital video recorder- (DVR-)like functions, such as play, fast forward, rewind, jump, pause, and stop. You can find performance bottlenecks by playing back the recorded data at any point in time.

For more information about PowerVP, visit:

http://pic.dhe.ibm.com/infocenter/powersys/v3r1m5/index.jsp?topic=%2Fp7ecul%2Fp7ecu_intro_powervp.htm.

AIX

The section introduces tools and techniques that are used for optimizing software for a combination of Power Systems and AIX. The intended audience for this section is software development teams. As such, this section does not address performance topics that are related to capacity planning, and system-level performance monitoring and tuning.

To download Java for AIX, visit the following website:

<http://www.ibm.com/developerworks/java/jdk/aix/>

For capacity planning, see the IBM Systems Workload Estimator, available here:

<http://www-912.ibm.com/estimator>

For system-level performance monitoring and tuning information for AIX, see *Performance management*, available here:

http://publib.boulder.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.prftungd/doc/prftungd/performance_management-kickoff.htm

The bedrock of any empirically based software optimization effort is a suite of repeatable benchmark tests. To be useful, such tests must be representative of the manner in which users interact with the software. For many commercial applications, a benchmark test simulates the actions of multiple users that drive a prescribed mix of application transactions. Here, the fundamental measure of performance is throughput (the number of transactions that are run over a period) with an acceptable response time. Other applications are more *batch-oriented*, where few jobs are started and the time that is taken to completion is measured. Whichever benchmark style is used, it must be repeatable. Within some small tolerance (typically a few percent), running the benchmark several times on the same setup yields the same result.

Tools and techniques that are employed in software performance analysis focus on pinpointing aspects of the software that inhibit performance. At a high level, the two most common inhibitors to application performance are:

- ▶ Areas of code that consume large amounts of CPU resources. This code is usually caused by using inefficient algorithms, poor coding practices, or inadequate compiler optimization
- ▶ Waiting for locks or external events. Locks are used to serialize execution through critical sections, that is, sections of code where the need for data consistency requires that only one software thread run at a time. An example of an external event is the system that is waiting for a disk I/O to complete. Although the amount of time that an application must wait for external events might be outside of the control of the application (for example, the time that is required for a disk I/O depends on the type of storage employed), simply being aware that the application is having to wait for such an event can open the door to potential optimizations.

CPU profiling

A CPU profiler is a performance tool that shows in which code CPU resources are being consumed. **tprof** is a powerful CPU profiler that encompasses a broad spectrum of profiling functionality:

- ▶ It can profile any program, library, or kernel extension that is compiled with C, C++, Fortran, or Java compilers. It can profile machine code that is created in real time by the JIT compiler.
- ▶ It can attribute time to processes, threads, subroutines (user mode, kernel mode, shared library, and Java methods), source statements, and even individual machine instructions.
- ▶ In most cases, no recompilation of object files is required.

Usage of **tprof** typically focuses on generating subroutine-level profiles to pinpoint code hotspots, and to examine the impact of an attempted code optimization. A common way to invoke **tprof** is as follows:

```
$ tprof -E -skeuz -x sleep 10
```

The **-E** flag instructs **tprof** to employ the performance monitoring unit (PMU) as the sampling mechanism to generate the profile. Using the PMU as the sampling mechanism provides a more accurate profile than the default time-based sampling mechanism, as the PMU sampling mechanism can accurately sample regions of kernel code where interrupts are disabled. The **s**, **k**, **e**, and **u** flags instruct **tprof** to generate subroutine-level profiles for shared library, kernel, kernel extension, and user-level activity. The **z** flag instructs **tprof** to report CPU time in the number of *ticks* (that is, samples), instead of percentages. The **-x sleep 10**

argument instructs **tprof** to collect profiling data during the running of the **sleep 10** command. This command collects profile data over the entire system (including all running processes) over a period of 10 seconds.

Excerpts from a **tprof** report are shown in Example B-1, Example B-2, and Example B-3.

Example B-1 is a breakdown of samples of the processes that are running on the system. When multiple processes have the same name, they have only one line in this report: the number of processes with that name is in the “Freq” column. “Total” is the total number of samples that are accumulated by the process, and “Kernel”, “User”, and “Shared” are the number of samples that are accumulated by the processes in kernel (including kernel extensions), user space, and shared libraries. “Other” is a catchall for samples that do not fall in the other categories. The most common scenario where samples wind up in “Other” is because of CPU resources that are being consumed by machine code that is generated in real time by the JIT compiler. The **-j** flag of **tprof** can be used to attribute these samples to Java methods.

Example: B-1 Excerpt from a tprof report - breakdown of samples of processes running on the system

Process	Freq	Total	Kernel	User	Shared	Other	
=====		=====	=====	=====	=====	=====	=====
wait		4	5810	5810	0	0	0
./version1	1	1672	35	1637	0	0	0
/usr/bin/tprof	2	15	13	0	2	0	0
/etc/syncd	1	2	2	0	0	0	0
/usr/bin/sh	2	2	2	0	0	0	0
swapper	1	1	1	0	0	0	0
/usr/bin/trcstop	1	1	1	0	0	0	0
rmcd	1	1	1	0	0	0	0
=====		=====	=====	=====	=====	=====	=====
Total		13	7504	5865	1637	2	0

Example B-2 is a breakdown of samples of the threads that are running on the system. In addition to the columns described in Example B-1 on page 208, this report has *PID* and *TID* columns that detail the process IDs and thread IDs.

Example: B-2 Excerpt from a tprof report - breakdown of threads that are running on the system

Process	PID	TID	Total	Kernel	User	Shared	Other	
=====	=====	=====	=====	=====	=====	=====	=====	=====
wait	16392	16393	1874	1874	0	0	0	0
wait	12294	12295	1873	1873	0	0	0	0
wait	20490	20491	1860	1860	0	0	0	0
./version1	245974	606263	1672	35	1637	0	0	0
wait	8196	8197	203	203	0	0	0	0
/usr/bin/tprof	291002	643291	13	13	0	0	0	0
/usr/bin/tprof	274580	610467	2	0	0	2	0	0
/etc/syncd	73824	110691	2	2	0	0	0	0
/usr/bin/sh	245974	606263	1	1	0	0	0	0
/usr/bin/sh	245976	606265	1	1	0	0	0	0
/usr/bin/trcstop	245976	606263	1	1	0	0	0	0
swapper	0	3	1	1	0	0	0	0
rmcd	155876	348337	1	1	0	0	0	0
=====	=====	=====	=====	=====	=====	=====	=====	=====
Total			7504	5865	1637	2	0	
Total Samples = 7504			Total Elapsed Time = 18.76s					

Example B-3 from the report gives the subroutine-level profile for the version1 program. In this simple example, all of the time is spent in `main()`.

Example: B-3 Excerpt from a tprof report - subroutine-level profile for the version1 program, with all time spent in main()

Profile: ./version1						
Total Ticks For All Processes (./version1) = 1637						
Subroutine	Ticks	%	Source	Address	Bytes	
=====	=====	=====	=====	=====	=====	=====
.main	1637	21.82	version1.c	350	536	

More information about using AIX **tprof** for Java programs is available in “Hot method or routine analysis” on page 223.

The functionality of **tprof** is rich. As such, it cannot be fully described in this guide. For complete **tprof** documentation, see *tprof Command*, available here:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.cmds/doc/aixcmds5/tprof.htm>

AIX trace-based analysis tools

Trace¹ is a powerful utility that is provided by AIX for collecting a time-sequenced log of operating system events on a Power Systems server. The AIX kernel and kernel extensions are richly instrumented with trace *hooks* that, when trace is activated, append trace records with context-relevant data, to a pinned, kernel-resident trace buffer. These records can be later read from that buffer and logged to a disk-resident file. Further utilities are provided to interpret and summarize trace logs and generate human-readable reports. The **tprof** CPU profiler is one such utility. Besides **tprof**, two of the most-commonly used trace-based utilities are **curt**² and **splat**.^{3,4}

The **curt** command takes as its input a trace collected using the AIX trace facility, and generates a report that breaks down how CPU time is consumed by various entities, including:

- ▶ Processes (grouped by process name)
- ▶ Individual processes
- ▶ Individual threads
- ▶ System calls (either on a system-wide or per-thread basis)
- ▶ Interrupts

One of the most useful reports from **curt** is the *System Calls Summary*. This report provides a system-wide summary of the system calls that are executed while the trace is collected. For each system call, the following information is provided:

- ▶ Count: The number of times the system call was run during the monitoring interval

¹ *trace Daemon*, available here:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.cmds/doc/aixcmds5/trace.htm>

² *CPU Utilization Reporting Tool (curt)*, available here:

http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.prftools/doc/prftools/idprftools_cpu.htm

³ *Simple performance lock analysis tool (splat)*, available here:

http://pic.dhe.ibm.com/infocenter/aix/v7r1/topic/com.ibm.aix.prftools/doc/prftools/idprftools_splat.htm

⁴ *splat Command*, available here:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.cmds/doc/aixcmds5/splat.htm>

- **Total Time:** Amount of CPU time (in milliseconds) consumed in running the system call
- **% sys time:** Percentage of overall CPU capacity that is spent in running the system call
- **Avg Time:** Average CPU time that is consumed for each execution of the system call
- **Min Time:** Minimum CPU time that is consumed during an execution of the system call
- **Max Time:** Maximum CPU time that is consumed during an execution of the system call
- **SVC:** Name and address of the system call

An excerpt from a System Calls Summary report is shown in Example B-4.

Example: B-4 System Calls Summary report (excerpt)

System Calls Summary						
Count	Total Time (msec)	% sys time	Avg Time (msec)	Min Time (msec)	Max Time (msec)	SVC (Address)
123647	3172.0694	14.60%	0.0257	0.0128	0.9064	kpread(2a2d5e8)
539	1354.6939	6.24%	2.5133	0.0163	4.1719	listio64(516ea40)
26496	757.6204	3.49%	0.0286	0.0162	0.0580	_esend(2a29f88)
26414	447.7029	2.06%	0.0169	0.0082	0.0426	_erecv(2a29e98)
9907	266.1382	1.23%	0.0269	0.0143	0.5350	kpwrite(2a2d588)
34282	167.8132	0.77%	0.0049	0.0032	0.0204	_thread_wait(2a28778)

As a first step, compare the mix of system calls to the expectation of how the application is expected to behave. Is the mix aligned with expectations? If not, first confirm that the trace is collected while the wanted workload runs. If the trace is collected at the correct time and the mix still differs from expectations, then investigate the application logic. Also, examine the list of system calls for potential optimizations. For example, if **select** or **poll** is used frequently, consider employing the pollset facility (see “pollset” on page 89).

As a further breakdown, **curt** provides a report of the system calls run by each thread. An example report is shown in Example B-5.

Example: B-5 system calls run by each thread

Report for Thread Id: 549305 (hex 861b9) Pid: 323930 (hex 4f15a)						
Process Name: procl						

Total Application Time (ms): 89.010297						
Total System Call Time (ms): 160.465531						
Total Hypervisor Call Time (ms): 18.303531						
Thread System Call Summary						
Count	Total Time (msec)	Avg Time (msec)	Min Time (msec)	Max Time (msec)	SVC (Address)	
492	157.0663	0.3192	0.0032	0.6596	listio64(516ea40)	
494	3.3656	0.0068	0.0002	0.0163	GetMultipleCompletionStatus(549a6a8)	
12	0.0238	0.0020	0.0017	0.0022	_thread_wait(2a28778)	
6	0.0060	0.0010	0.0007	0.0014	thread_unlock(2a28838)	
4	0.0028	0.0007	0.0005	0.0008	thread_post(2a288f8)	

Another useful report that is provided by **curt** is the *Pending System Calls Summary*. This summary shows the list of threads that are in an unfinished system call at the end of the trace. An example report is given in Example B-6.

Example: B-6 Threads that are in an unfinished system call at the end of the trace

Pending System Calls Summary

Accumulated Time (msec)	SVC (Address)	Procname (Pid Tid)
0.0082	GetMultipleCompletionStatus(549a6a8)	proc1(323930 532813)
0.0089	_nsleep(2a28d30)	proc2(270398 545277)
0.0054	_thread_wait(2a28778)	proc1(323930 549305)
0.0088	GetMultipleCompletionStatus(549a6a8)	proc1(323930 561437)
3.3981	listio64(516ea40)	proc1(323930 577917)
0.0130	kpwrite(2a2d588)	proc1(323930 794729)

For each thread in an unfinished system call, the following items are provided:

- ▶ The accumulated time in the system call
- ▶ The name of the system call (followed by the system call address in parentheses)
- ▶ The process name, followed by the Process ID and Thread ID in parentheses

This report is useful in determining what system calls are blocking threads from proceeding. For example, threads appearing in this report with an unfinished **recv** call are waiting on data to be received over a socket.

Another useful trace-based tool is **splat**, which is the Simple Performance Lock Analysis Tool. The **splat** tool provides reports about the usage of kernel and application (pthread-level) locks. At the pthread level, **splat** can report about the usage of pthread synchronizers: mutexes, read/write locks, and condition variables. Importantly, **splat** provides data about the degree of contention and blocking on these objects, an important consideration in creating highly scalable and pthread-based applications.

The pthread library instrumentation does not provide names or classes of synchronizers, so the addresses are the only way that you have to identify them. Under certain conditions, the instrumentation can capture the return addresses of the function call stack, and these addresses are used with the output of the **gensyms** tool to identify the call chains when these synchronizers are created. The creation and deletion times of the synchronizer can sometimes be determined as well, along with the ID of the pthread that created them.

An example of a mutex report from **splat** is shown in Example B-7.

Example: B-7 Mutex report from splat

[pthread MUTEX] ADDRESS: 00000000F0154CD0

Parent Thread: 0000000000000001 creation time: 26.232305

Pid: 18396 Process Name: trcstop

Creation call-chain =====

```

00000000D268606C      .pthread_mutex_lock
00000000D268EB88      .pthread_once
00000000D01FE588      .__libs_init
00000000D01EB2FC      dne_callbacks
00000000D01EB280      .libc_declare_data_functions
00000000D269F960      .pth_init_libc
00000000D268A2B4      .pthread_init
00000000D01EAC08      .__modinit
000000001000014C      .__start
  
```

Acqui- sitions	Miss Rate	Spin Count	Wait Count	Busy Count	Secs Held		Percent Held (26.235284s)			
					CPU	Elapsed	Real CPU	Real Elapsed	Comb Spin	Real Wait
1	0.000	0	0	0	0.000006	0.000006	0.00	0.00	0.00	0.00

Depth	Min	Max	Avg
-------	-----	-----	-----

SpinQ	0	0	0									
WaitQ	0	0	0									
Recursion	0	1	0									
	Acqui-	Miss	Spin	Wait	Busy	Percent Held of Total Time						
PThreadID	sitions	Rate	Count	Count	Count	CPU	Elapse	Spin	Wait			
~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~			
1	1	0.00	0	0	0	0.00	0.00	0.00	0.00			
	Acqui-	Miss	Spin	Wait	Busy	Percent Held of Total Time						
Function Name	sitions	Rate	Count	Count	Count	CPU	Elapse	Spin	Wait	Return Address	Start Address	
Offset	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	
~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	~~~~~	
.pthread_once	0	0.00	0	0	0	99.99	99.99	0.00	0.00	00000000D268EC98	00000000D2684180	
.pthread_once	1	0.00	0	0	0	0.01	0.01	0.00	0.00	00000000D268EB88	00000000D2684180	

In addition to the common header information and the [pthread_mutex_t] identifier, this report lists the following lock details:

Parent thread	Pthread ID of the parent pthread	
Creation time	Elapsed time in seconds after the first event recorded in trace (if available)	
Deletion time	Elapsed time in seconds after the first event recorded in trace (if available)	
PID	Process identifier	
Process Name	Name of the process using the lock	
Call-chain	Stack of called methods (if available)	
Acquisitions	The number of times the lock was acquired in the analysis interval	
Miss Rate	The percentage of attempts that failed to acquire the lock	
Spin Count	The number of unsuccessful attempts to acquire the lock	
Wait Count	The number of times a thread is forced into a suspended wait state while waiting for the lock to come available	
Busy Count	The number of trylock calls that returned busy	
Seconds Held	This field contains the following subfields:	
	CPU	The total number of processor seconds the lock is held by a running thread.
	Elapse(d)	The total number of elapsed seconds the lock is held, whether the thread was running or suspended.
Percent Held	This field contains the following subfields:	
	Real CPU	The percentage of the cumulative processor time the lock was held by a running thread.
	Real Elapsed	The percentage of the elapsed real time the lock is held by any thread, either running or suspended.
	Comb(ined) Spin	The percentage of the cumulative processor time that running threads spend spinning while it tries to acquire this lock.
	Real Wait	The percentage of elapsed real time that any thread was waiting to acquire this lock. If two or more threads are waiting simultaneously, this wait time is only charged one time. To learn how many threads are waiting simultaneously, look at the WaitQ Depth statistics.

Depth	This field contains the following subfields:	
	SpinQ	The minimum, maximum, and average number of threads that are spinning on the lock, whether running or suspended, across the analysis interval
	WaitQ	The minimum, maximum, and average number of threads that are waiting on the lock, across the analysis interval
	Recursion	The minimum, maximum, and average recursion depth to which each thread held the lock

Finding alignment issues

Improperly aligned code or data can cause performance degradation. By default, the IBM compilers and linkers correctly align code and data, including stack and statically allocated variables. Incorrect typecasting can result in references to storage that are not correctly aligned. There are two types of alignment issues to be concerned with:

- ▶ Alignment issues that are handled by microcode in the POWER7 processor
- ▶ Alignment issues that are handled through alignment interrupts.

Examples of alignment issues that are handled by microcode with a performance penalty in the POWER7 processor are loads that cross a 128-byte boundary and stores that cross a 4 KB page boundary. To give an indication of the penalty for this type of misalignment, on a 4 GHz processor, a nine-instruction loop that contains an 8 byte load that crosses a 128-byte boundary takes double the time of the same loop with the load correctly aligned.

Alignment issues that are handled by microcode can be detected by running **hpmcount** or **hpmstat**. The **hpmcount** command is a command-line utility that runs a command and collects statistics from the POWER7 PMU while the command runs. To detect alignment issues that are handled by microcode, run **hpmcount** to collect data for group 38. An example is provided in Example B-8.

Example: B-8 Example of the results of the hpmcount command

# hpmcount -g 38 ./unaligned		
Group: 38		
Counting mode: user		
Counting duration: 21.048874056 seconds		
PM_LSU_FLUSH_ULD (LRQ unaligned load flushes)	:	4320840034
PM_LSU_FLUSH_UST (SRQ unaligned store flushes)	:	0
PM_LSU_FLUSH_LRQ (LRQ flushes)	:	450842085
PM_LSU_FLUSH_SRQ (SRQ flushes)	:	149
PM_RUN_INST_CMPL (Run instructions completed)	:	19327363517
PM_RUN_CYC (Run cycles)	:	84219113069
Normalization base: time		
Counting mode: user		
Derived metric group: General		
[] Run cycles per run instruction	:	4.358

The **hpmstat** command is similar to **hpmcount**, except that it collects performance data on a system-wide basis, rather than just for the execution of a command.

Generally, scenarios in which the ratio of (*LRQ unaligned load flushes + SRQ unaligned store flushes*) divided by *Run instructions completed* is greater than 0.5% must be further investigated. The **tprof** command can be used to further pinpoint where in the code the unaligned storage references are occurring. To pinpoint unaligned loads, the

-E **PM_MRK_LSU_FLUSH_ULD** flag is added to the **tprof** command line, and to pinpoint unaligned stores, the -E **PM_MRK_LSU_FLUSH_UST** flag is added. When these flags are used, **tprof** generates a profile where unaligned loads and stores are sampled instead of time-based sampling.

Examples of alignment issues that cause an alignment interrupt include execution of a **lmmw** or **lwarx** instruction on a non-word-aligned boundary. These issues can be detected by running **alstat**. This command can be invoked with an interval, which is the number of seconds between each report. An example is presented in Example B-9.

Example: B-9 Alignment issues can be addressed with the alstat command

```
> alstat 5
Alignment  Alignment
SinceBoot   Delta
    2016      0
    2016      0
    2016      0
    2016      0
    2016      0
    2016      0
```

The key metric in the **alstat** report is the Alignment Delta. This metric is the number of alignment interrupts that occurred during the interval. Non-zero counts in this column merit further investigation with **tprof**. Invoking **tprof** with the -E **ALIGNMENT** flag generates a profile that shows where the unaligned references are occurring.

For more information, see *alstat Command*, available here:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.commands/doc/aixcmds1/alstat.htm>

Finding emulation issues

Over the 20+ year evolution of the Power instruction set, a few instructions were removed. Instead of trapping programs that run these instructions, AIX emulates them in the kernel, although with a significant processing impact. Generally, programs that are written in a third-generation language (for example, C and C++) and compiled with an up-to-date compiler do not contain these emulated instructions. However, older binary files or older hand-written assembly language might contain such instructions, and because they are silently emulated by AIX, the performance penalty might not be readily apparent.

The **emstat** command detects the presence of these instructions. Like **alstat**, it is invoked with an interval, which is the number of seconds between reports. An example is shown in Example B-10.

Example: B-10 The emstat command detects the presence of emulated instructions

```
> emstat 5
Emulation  Emulation
SinceBoot   Delta
    0        0
    0        0
    0        0
    0        0
    0        0
```

The key metric is the Emulation Delta (the number of instructions that are emulated during each interval). Non-zero values merit further investigation. Invoking **tprof** with the **-E EMULATION** flag generates a profile that shows where the emulated instructions are.

For more information, see *emstat Command*, available here:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.cmds/doc/aixcmds2/emstat.htm>

hpmstat, hpmcount, and tprof -E

The POWER7 processor provides a powerful on-chip PMU that can be used to count the number of occurrences of performance-critical processor events. A rich set of events is countable; examples include level 2 and level 3 d-cache misses, and cache reloads from local, remote, and distant memory. *Local memory* is memory that is attached to the same POWER7 processor chip that the software thread is running on. *Remote memory* is memory that is attached to a different POWER7 processor that is in the same central electronic complex (CEC) (that is, the same node or building block in the case of a multi-CEC system, such as a Power 780) that the software thread is running on. *Distant memory* is memory that is attached to a POWER7 processor that is in a different CEC from the CEC the software thread is running on.

Two commands exist to count PMU events: **hpmcount** and **hpmstat**. The **hpmcount** command is a command-line utility that runs a command and collects statistics from the PMU while the command runs. The **hpmstat** command is similar to **hpmcount**, except that it collects performance data on a system-wide basis, rather than just for the execution of a command.

Further documentation about **hpmcount** and **hpmstat** can be found here:

- ▶ <http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.cmds/doc/aixcmds2/hpmcount.htm>
- ▶ <http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.cmds/doc/aixcmds2/hpmstat.htm>

In addition to simply counting processor events, the PMU can be configured to sample instructions based on processor events. With this capability, profiles can be generated that show which parts of an application are experiencing specified processor events. For example, you can show which subroutines of an application are generating level 2 or level 3 cache misses. The **tprof** profiler includes this functionality through the **-E** flag, which allows a PMU event name to be provided to **tprof** as the sampled event. The list of PMU events can be generated by running **pm1ist -c -1**. Whenever possible, perform profiling using *marked* events, as profiling using marked events is more accurate than using unmarked events. The marked events begin with the prefix **PM_MRK_**.

For more information about using the **-E** flag of **tprof**, go to this website:

<http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.cmds/doc/aixcmds5/tprof.htm>

Linux

The section introduces tools and techniques used for optimizing software on the combination of Power Systems and Linux. The intended audience for this section is software development teams.

To download Java for Linux, visit the following website:

For Linux: <http://www.ibm.com/developerworks/java/jdk/linux/>

Empirical performance analysis using the IBM software development kit (SDK) for PowerLinux

After you apply the best high-level optimization techniques, a deeper level of analysis might be required to gain more performance improvements. You can use the IBM SDK for PowerLinux to help you gain these improvements.

The IBM SDK for PowerLinux is a set of tools that support:

- ▶ Hot spot analysis
- ▶ Analysis of ported code for missed platform-specific optimization
- ▶ Whole program analysis for coding issues, for example, pipeline hazards, inlining opportunities, early exits and hidden path length, devirtualization, and branch prediction hints
- ▶ Lock contention and IO delay analysis

The *IBM SDK for PowerLinux* can be found here:

<http://www.ibm.com/support/customercare/sas/f/lopdiags/sdklop.html>

The SDK provides an Eclipse C/C++ IDE with Linux tools integrations. The SDK provides graphical presentation and source code view integration with Linux execution profiling (**gprof/OProfile/Perf**), malloc and memory usage (**valgrind**), pthread synchronization (**helgrind**), SystemTap tapsets, and tapset development.

Hotspot analysis

We suggest that you profile the application and look for hotspots. When you run the application under one or more representative workloads, use a hardware-based profiling tool such as **OProfile**. **OProfile** can be run directly as a command-line tool or under the IBM SDK for PowerLinux.

The **OProfile** tools can monitor the whole system (LPAR), including all the tasks and the kernel. This action requires root authority, but is the best way to profile the kernel and complex applications with multiple cooperating processes. **OProfile** is fully enabled to take samples using the full set of the PMU events (run **ophelp** for a complete list of events). **OProfile** can produce text file reports organized by process, program and libraries, function symbols, and annotated source file and line number or machine code disassembly.

The IBM SDK for PowerLinux can profile applications that are associated with Eclipse projects. The SDK automates the setup and running of the profile, but is restricted to a single application, its libraries, and direct kernel calls. The SDK is easier to use, as it is hierarchically organized by percentage with program, function symbol, and line number. Clicking the line number in the profile pane *jumps* the source view pane to the matching source file and line number. This action simplifies edit, compile, and profile tuning activities.

The whole system profile is a good place to start. You might find that your application is consuming most of the CPU cycles, and deeper analysis of the application is the next logical step. The IBM SDK for PowerLinux provides a number of helpful tools, including integrated application profiling (**OProfile** and **valgrind**), Migration Assistant, and the Source Code Advisor.

High kernel usage

If the bulk of the CPU cycles are consumed in the kernel or runtime libraries that are not part of your application, then a different type of analysis is required. If the kernel is consuming significant cycles, then the application might be I/O or lock contention bound. This situation can occur when an application moves to larger systems (higher core count) and fails to scale up.

I/O bound applications can be constrained by small buffer sizes or a poor choice of an access method. One issue to look for is applications that use local loopback sockets for interprocess communications (IPC). This situation is common for applications that are migrating from early scale-out designs to larger systems (and core-count). The first application change is to choose a lighter weight form of IPC for in-system communications.

Excessive locking or poor lock granularity can also result in high kernel usage (in the kernel's `spin_lock`, `futex`, and scheduler components) when applications move to larger system configurations. This situation might require adjusting the application lock strategy and possibly the type of lock mechanism that is used as well:

- ▶ POSIX `pthread_mutex` and `pthread_rwlock` locks are complex and heavy, and POSIX semaphores are simpler and lighter.
- ▶ Use `trylock` forms to spin in user mode for a limited time when appropriate. Use this technique when there is normally a finite lock hold time and limited contention for the resource. This situation avoids context switch and scheduler impact in the kernel.
- ▶ Reserve POSIX `pthread_spinlock` and `sched_yield` for applications that have exclusive use of the system and with carefully designed thread affinity (assigning specific threads to specific cores).
- ▶ The compiler provides inline functions (`__sync_fetch_and_add`, `__sync_fetch_and_or`, and so on) that are better suited for simple atomic updates than POSIX lock and unlock. Use thread local storage, where appropriate, to avoid locking for thread safe code.

Using the IBM SDK for PowerLinux Trace Analyzer

The IBM SDK for PowerLinux provides tools, including the SystemTap and pthread monitor, for tracking I/O and lock usage of a running application. The higher level Trace Analyzer tools can target a specific application for combined SystemTap syscall trace and Lock Trace. The resulting trace information is correlated for time strip display and analysis within the tool.

High library usage

If libraries are consuming significant cycles, then you must determine these possibilities:

- ▶ Those libraries are part of your application, provided by a third party, or the Linux distribution.
- ▶ There are alternative libraries that are better optimized.
- ▶ You can recompile those libraries at a higher optimization.

Libraries that are part of your application require the same level of empirical analysis as the rest of your application (by using source profiling and the Source Code Advisor (SCA)). Libraries that are used by but are not part of your application imply a number of options and strategies:

- ▶ Most open source packages in the Linux environment are compiled with optimization level `-O2` and tend to avoid additional (higher level GCC) compiler options. This configuration might be sufficient for a CISC processor with limited register resources, but not sufficient for a RISC based register-rich processor, such as POWER7 and POWER8.

- ▶ A RISC-based, superscalar, out-of-order execution processor chip such as POWER8 and POWER8 requires more aggressive inlining and loop-unrolling to capitalize on the larger register set and superscalar design point. Also, automatic vectorization is not enabled at this lower (**-O2**) optimization level, and so the vector registers and ISA feature go unused.
- ▶ In GCC, you must specify the **-O3** optimization level and inform the compiler that you are running on a newer processor chip with the Vector ISA extensions. In fact, with GCC, you need both **-O3** and **-mcpu=power7** for the compiler to generate code that capitalizes on the new VSX feature of POWER7. You will need both **-O3** and **-mcpu=power8** for the compiler to take advantage of the latest VSX instructions implemented on POWER8.

One source of optimized libraries is the IBM Advance Toolchain for PowerLinux. The Advance Toolchain provides alternative runtime libraries for all the common POSIX C language, Math, and pthread libraries that are highly optimized (**-O3** and **-mcpu=**) for multiple Power platforms (including POWER7 and POWER8). The Advance Toolchain run time RPM provides multiple CPU tuned library instances and automatically selects the specific library version that is optimized for the specific POWER5, POWER6, POWER7, or POWER8 machine.

If there are specific open source or third-party libraries that are dominating the execution profile of your application, you must ask the distribution or library product owner to provide a build using higher optimization. Alternatively, for open source library packages, you can build your own optimized binary version of those packages.

Deeper empirical analysis

If simple recompilation with higher optimization options or even a more capable compiler does not provide acceptable performance, then deeper analysis is required. The IBM SDK for PowerLinux integrates the following analysis tools:

- ▶ Migration Assistant analysis, non-performing codes, and data types
- ▶ Application-specific hotspot profiling
- ▶ SCA analysis for non-performing code idioms and induced execution hazards

The Migration Assistant analyzes the source code directly and does not require a running binary application for analysis. Profiling and the SCA do require compiled application binary files and an application-specific benchmark or repeatable workload for analysis.

The Migration Assistant

For applications that originate on another platform, the Migration Assistant (MA) can identify non-portable code that must be addressed for a successful port to Power Systems. The MA uses the Eclipse infrastructure to analyze:

- ▶ Data endian dependent unions and structures
- ▶ Casts with potential endian issues
- ▶ Non-portable data types
- ▶ Non-portable inline assembler code
- ▶ Non-portable or arch dependent compiler built-ins
- ▶ Proprietary or architectural-specific APIs

Program usage of non-portable data types and an inline assembler can cause poor performance on the POWER processor, which always must be investigated and addressed.

For example, the long double data type is supported for both Intel x86 and Power, but has a different size, data range, and implementation. The x86 80-bit Floating Point format is implemented in hardware and is usually faster than (although not compatible with) the AIX long double, which is implemented as an algorithm using two, 64-bit doubles. Neither one is fully IEEE-compliant, and both must be avoided in cross-platform application codes and libraries.

Another example is small Intel specific optimization using inline x86 assembler and conditionally providing a generic C implementation for other platforms. In most cases, GCC provides an equivalent built-in function that generates the optimal code for each platform. Replacing inline assembler with GCC built-in functions makes the application more portable and provides equivalent or better performance on all platforms.

To use the MA tool, complete the following steps:

1. Import your project into the SDK.
2. Select Project **properties**.
3. Check the **Linux/x86 to PowerLinux application Migration** check box under C/C++ General/Code Analysis.
4. Right click the project name, and select **Run Migration Advisor**.

Hotspot profiling

IBM SDK for PowerLinux integrates the Linux **OProfile** hardware event profiling with the application source code view. This configuration is a convenient way to do hotspot analysis. The integrated Linux Tools profiler focuses on an application that is selected from the current SDK project.

After you run the application, the SDK opens an **OProfile** tab in console window. This window shows a nested set of *twisties*, starting with the *event* (cycles by default), then *program/library*, *function*, and *source line* (within function). The developer drills-down by opening the twisties in the profile window, opening the next level of detail. Items are ordered by profile frequency with highest frequency first. Clicking the function or line number entries in the profile window causes the source view to *jump* to the corresponding source file or line number.

This process is a convenient way to do hotspot analysis, focusing only on the top three to five items at each level in the profile. Examine the source code for algorithmic problems, excess conversions, unneeded debug code, and so on, and make the appropriate source code changes.

With your application code (or subset) imported in to the SDK, it is easy to edit, compile, and profile code changes and verify improvements. As the developer makes code improvements, the hotspots in the profile change. Repeat this process until performance is satisfactory or all the profile entries at the function level are in the low single digits.

To use the integrated profiler, right-click the project and select **Profile As → Profile with OProfile**. If you project contains multiple applications or the application needs setup or inputs to run the specific workload, then create Profile Configurations as needed.

Detailed analysis with the Source Code Advisor (SCA)

Hotspot analysis might not find all of the latent performance problems, especially coding style and some machine-specific hazards. These problems tend to be diffused across the application, and do not show up in hotspot analysis. Common examples of machine hazards include address translation, cache misses, and branch miss-predictions.

Complex C++ applications or C programs that use object-based techniques might see performance issues that are related to using many small functions of indirect calls. Unless the compiler or optimizer can see the whole program or library, it cannot prove that it is safe to optimize these cases. However, it is possible for the developer to manually optimize at the source level, as the developer knows the original intent or actual usage in context.

The SCA can find and suggest solutions for many of these coding style and machine hazards. The process generates a journal that associates performance problems (including hazards) with specific source file and line numbers.

The SCA window has a drill-down hierarchy similar to the profile window described in “Hotspot profiling” on page 219. The SCA window is organized as a list of problem categories, and then nested twisties, for affected functions and source line numbers within functions. Functions and lines are ordered by the percent of overall contribution to execution time. Associated with each problem is a plain language description and suggested solution that describes a source change or compiler or linker options that are expected to resolve the problem. Clicking the line number item *jumps* the source display to the associated source file and line number for editing.

SCA uses the Feedback Directed Program Restructuring (FDPR) tool to instrument your application (or library) for code and data flow trace when you run a workload. The resulting FDPR journal is used to drive the SCA analysis. Running FDPR and retrieving the journal is automated by clicking **Profile as** → **Profile with Source Code Advisor**.

Pipeline stall analysis with the cycles per instruction (CPI) breakdown tool

The cycles per instruction (CPI) metric is a measure of the average processor clock cycles that are needed to complete an instruction. The CPI value is a measure of processor performance and, in a modern processor such as the POWER processor, a high value can indicate poor performance due to a high ratio of stalls in the execution pipeline. By collecting information from the processor’s PMU, those events and derived metrics can be mapped to the CPU functional units (for example, branch, load or store, or floating point), where they occurred. These events and metrics can be represented in a hierarchical breakdown of cycles, called the CPI breakdown model (CBM). Further information about CPI metric and pipeline analysis is located in *Commonly Used Metrics for Performance Analysis*, available here:

<https://www.power.org/documentation/commonly-used-metrics-for-performance-analysis/> (registration required)

The IBM SDK for PowerLinux delivers with the CPI breakdown tool for automating the collection of PMU stall events and for building a CBM representation of application execution. After you run the application, the CPI breakdown tool opens a CBM view in the default Eclipse perspective. This view shows a breakdown of stall events and metrics, along with their contribution percentage and description. A sample is shown in Figure B-6.

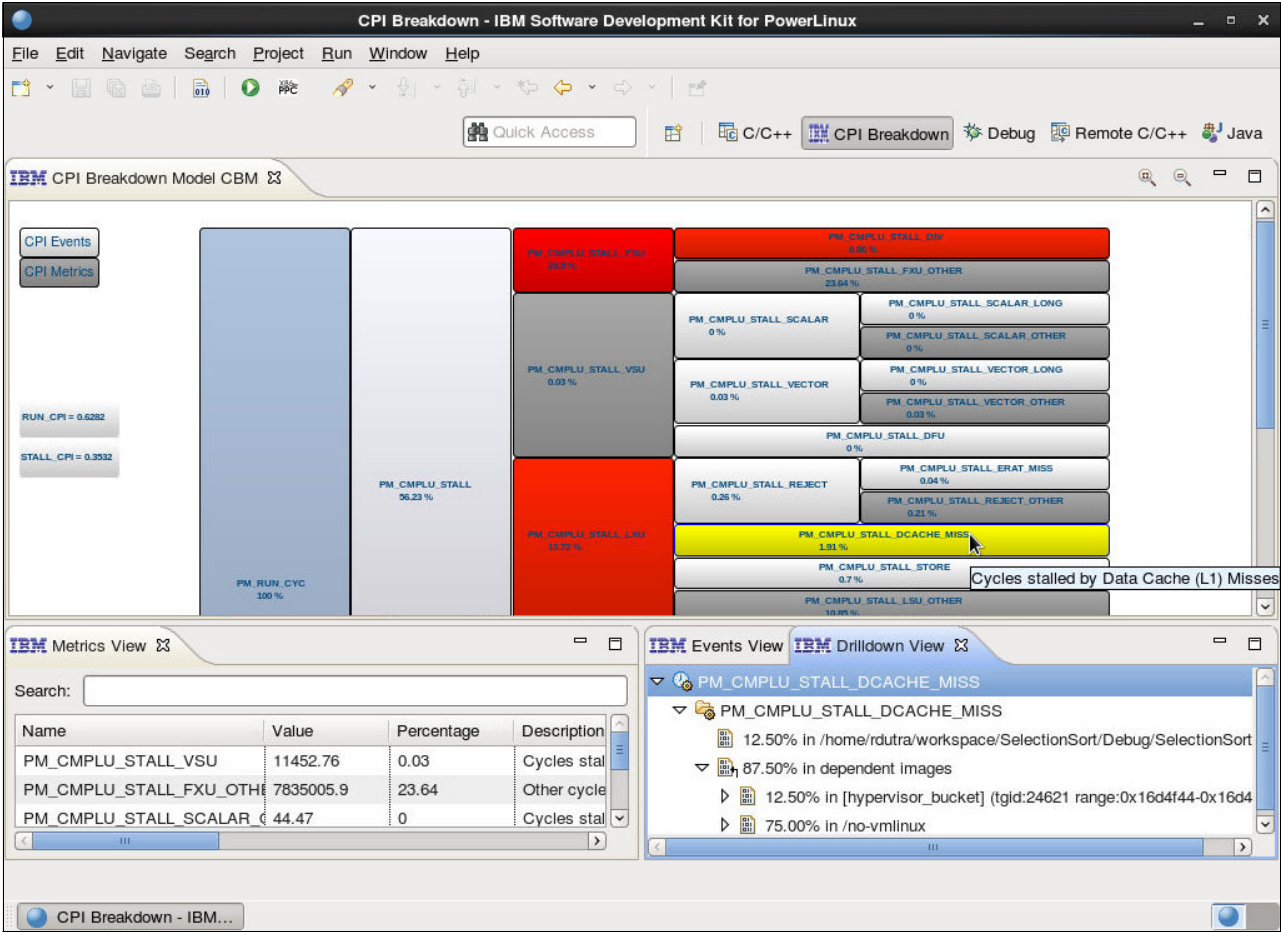


Figure B-6 The CPI Breakdown tool perspective

In the CBM view (see Figure B-6), click in any of the squares to open the drill-down that shows a nested set of twisties, including the event, program or library, function, and source line. This can be seen in Figure B-6, bottom-right. As you drill down, the items are ordered by profile frequency, with highest frequency first. Click a function or line number in the profile window to open the source view and jump to the corresponding source file and line number.

Use the CPI Breakdown tool to measure application behavior in the Power processor and for hotspot analysis. The tool assists in finding the CPU functional units with a high ratio of stalls and the corresponding chunks of source code that are likely to be the cause of performance degradation of your application.

Java (either AIX or Linux)

Focused empirical analysis of Java applications involves gathering specific types of performance information, making and assessing changes, and repeating the process. The specific areas to consider, the types of performance information to gather, and the tools to use, are described in this section.

To download Java for AIX and Linux, visit the following websites:

For AIX: <http://www.ibm.com/developerworks/java/jdk/aix/>

For Linux: <http://www.ibm.com/developerworks/java/jdk/linux/>

32-bit or 64-bit JDK

All other things being equal, a 64-bit JDK using **-Xcompressedrefs** generally has about 5% lower performance than does a 32-bit JDK. Without the **-Xcompressedrefs** option, a 64-bit JDK might have 10% or more reduced performance, which is compared to a 32-bit JDK. Give careful consideration to the choice of a 32-bit or 64-bit JVM. It is *not* a good choice to take an application that suffers from excessive object allocation rates and switch to a 64-bit JVM simply to allow a larger heap size. The references in the related tools and analysis techniques information can be used to diagnose object allocation issues in an application.

For more information about this topic, see these sections:

- ▶ “Verbose GC Log” on page 222
- ▶ 8.2, “32-bit versus 64-bit Java” on page 162.

Java heap size, and garbage collection (GC) policies and parameters

The performance of Java applications is often influenced by the heap size, GC policy, and GC parameters. Try different combinations, which are guided by appropriate data gathering and analysis. Various tools and diagnostic options are available that can provide detailed information about the state of the JVM. The information that is provided can be used to guide tuning decisions to maximize performance for an application or workload.

Verbose GC Log

The verbose GC log is a key tool to understanding the memory characteristics of a particular workload. The information that is provided in the log can be used to guide tuning decisions to minimize GC impact and improve overall performance. Logging can be activated with the **-verbose:gc** option and is directed to the command terminal. Logging can be redirected to a file with the **-Xverbosegclog:<file>** option.

Verbose logs capture many types of GC events, such as regular GC cycles, allocation failures, heap expansion and contraction, events related to concurrent marking, and scavenger collections. Verbose logs also show the approximate length of time many events take, the number of bytes processed (if applicable), and other relevant metrics. Information relevant to many of the tuning issues for GC can be obtained from the log, such as appropriate GC policies, optimal constant heap size, optimal min and max free space factors, and growth and shrink sizes. For a detailed description of verbose log output, consult the material on this subject in the *Diagnostics Guide for IBM SDK, Java Technology Edition, Version 6*, available here:

<http://publib.boulder.ibm.com/infocenter/javasdk/v6r0/topic/com.ibm.java.doc.diagnostics.60/homepage/plugin-homepage-java6.html>

Garbage collection (GC) and memory visualizer

For large, long-running workloads, verbose logs can quickly grow in size, making them difficult to work with and to analyze an application's behavior over time. The GC and memory visualizer is a tool that can parse verbose GC logs and present them in a visual manner using graphs and other diagrams, allowing trends and totals to be easily and quickly recognized. The graphs can be used to determine the minimum and maximum heap usage, growth and shrink rates over time, and identify oscillating behaviors. This information can be especially helpful when you choose optimal GC parameters. The GC and memory visualizer can also compare multiple logs side by side, which can aid in testing various options in isolation and determining their effects.

For more information about the GC and memory visualizer, see *IBM Monitoring and Diagnostic Tools for Java - Garbage Collection and Memory Visualizer Version 2.7*, available here:

<http://www.ibm.com/developerworks/java/jdk/tools/gcmv/>

Java Health Center

The Java Health Center is the successor to both the GC and memory visualizer and the Java Lock Monitor. It is an all-in-one tool that provides information about GC activity, memory usage, and lock contention. The Health Center also functions as a profiler, providing sample-based statistics on method execution. The Health Center functions as an agent of the JVM being monitored and can provide information throughout the life of a running application.

For more information about the Java Health Center, see *Java diagnostics, IBM style, Part 5: Optimizing your application with the Health Center*, available here:

<https://www.ibm.com/developerworks/java/library/j-ibmtools5>

For more information, see 8.4, “Java garbage collection tuning” on page 168.

Hot method or routine analysis

A CPU profile shows a breakdown of the time that is spent in Java methods and JNI or system routines. Investigate any hot methods or routines to determine if the concentration of execution time in them is warranted or whether there is poor coding or other issues.

Some tools and techniques for this analysis include:

- ▶ AIX **tprof** profiling. For more information, see *tprof Command*, available here:
<http://publib.boulder.ibm.com/infocenter/aix/v7r1/index.jsp?topic=/com.ibm.aix.cmds/doc/aixcmds5/tprof.htm>
- ▶ Linux **OProfile** profiling. For more information about **OProfile**, see the following resources:
 - *Getting started with OProfile on PowerLinux* (resource page), available here:
<http://pic.dhe.ibm.com/infocenter/lxinfo/v3r0m0/index.jsp?topic=%2Fliacf%2Foprofgetstart.htm>
 - *Getting started with OProfile on PowerLinux*, available here:
http://pic.dhe.ibm.com/infocenter/lxinfo/v3r0m0/topic/liacf/oprofile_pdf.pdf
 - *OProfile results with JIT samples*, available here:
<http://oprofile.sourceforge.net/doc/getting-jit-reports.html>

- *Java Performance on POWER7*, available here:
https://www.ibm.com/developerworks/community/wikis/home?lang=en#!/wiki/W51a7ffc4dfd_4b40_9d82_446ebc23c550/page/Java%20Performance%20on%20POWER7
- *OProfile manual*, available here:
<http://oprofile.sourceforge.net/doc/index.html>

General information about running the profiler and interpreting the results are contained in the sections on profiling in “AIX” on page 206 and “Linux” on page 215. For Java profiling, additional Java options are required to be able to profile the machine code that is generated for methods by the JIT compiler:

- ▶ AIX 32-bit: **-agentlib:jpa=instructions=1**
- ▶ AIX 64-bit: **-agentlib:jpa64=instructions=1**
- ▶ Linux OProfile: **-agentlib:jvmti_oprofile**

The entire execution of a Java program can be profiled, for example on AIX by running the following command:

```
tprof -ujeskl -A -I -E -x java ...
```

However, it is more common to profile Java after a warm-up period so that JIT compilation activity has generally completed. To profile after a warm-up, start Java and wait an appropriate interval until steady-state performance is reached, which is anywhere from a few seconds to a few minutes for large applications. Then, invoke the profiler, for example, on AIX, by running the following command:

```
tprof -ujeskl -A -I -E -x sleep 60
```

On Linux, **OProfile** can be used in a similar fashion; for more information, see “Java profiling example”, and follow the appropriate documentation in the resources included in this section.

Java profiling example

Example B-11 contains a sample Java program that is profiled on AIX and Linux. This program does some meaningless work and is purposely poorly written to illustrate lock contention and GC impact in the profile. The program creates three threads but serializes their execution by having them attempt to lock the same object. One thread at a time acquires the lock, forcing the other two threads to wait until they can get the lock and run the code that is protected by the synchronized statement in the `doWork` method. While they wait to acquire the lock, the threads initially use *spin locking*, repeatedly checking if the lock is free. After a suitable amount of spinning, the threads block rather than continuing to use CPU resources.

Example: B-11 Sample Java program

```
public class ProfileTest extends Thread {

    static Object o; /* used for locking to serialize threads */
    static Double A[], B[], C[];
    static int Num=1000;

    public static void main(String[] args) {
        o = new Object();
        new ProfileTest().start(); /* start 3 threads */
        new ProfileTest().start(); /* each thread executes the "run" method */
        new ProfileTest().start();
    }
}
```

```

public void run() {
    double sum = 0.0;
    for (int i = 0; i < 50000; i++) {
        sum += doWork(); /* repeatedly do some work */
    }
    System.out.println("sum: "+sum); /* use the results of the work */
}

public double doWork() {
    double d;
    synchronized (o) { /* serialize the threads to create lock contention */
        A = new Double [Num];
        B = new Double [Num];
        C = new Double [Num];
        initialize();
        calculate();
        d = C[0].doubleValue();
    }
    return(d); /* use the calculated values */
}

public static void initialize() {
    /* Initialize A and B. */
    for (int i = 0; i < Num; i++) {
        A[i] = new Double(Math.random()); /* use new to create objects */
        B[i] = new Double(Math.random()); /* to force garbage collection */
    }
}

public static void calculate() {
    for (int i = 0; i < Num; i++) {
        C[i] = new Double(A[i].doubleValue() * B[i].doubleValue());
    }
}
}

```

The program also uses the Double class, creating many short-lived objects by using **new**. By running the program with a small Java heap, GC frequently is required to free the Java heap space that is taken by the Double objects that are no longer in use.

Example B-12 shows how this program was run and profiled on AIX. 64-bit Java was used with the options **-Xms10m** and **-Xmx10m** to specify the size of the Java heap. The profile that is generated appears in the java.prof file.

Example: B-12 Results of running tprof on AIX

```
# tprof -ujeskl -A -I -E -x java -Xms10m -Xmx10m -agentlib:jpa64=instructions=1 ProfileTest
```

Starting Command java -Xms10m -Xmx10m -agentlib:jpa64=instructions=1 ProfileTest

```

sum: 12518.481782746869
sum: 12507.63528674597
sum: 12526.320955364286
stopping trace collection.
Sun Oct 30 15:04:21 2011
System: AIX 6.1 Node: el9-90-28 Machine: 00F603F74C00

```

Generating java.trc
Generating java.syms

Generating java.prof

Example B-13 and Example B-14 contain excerpts from the java.prof file that is created on AIX. The notable elements of the profile are:

- **Lock contention impact:** The impact of spin locking is shown in Example B-13 as ticks in the libj9jit24.so helper routine jitMonitorEntry, in the AIX pthreads library libpthreads.a, and in the AIX kernel routine **_check_lock**. This Java program clearly has excessive lock contention with jitMonitorEntry consuming 26.66% of the ticks in the profile. jitMonitorEntry and other routines, such as jitMethodMonitorEntry, indicate spin locking at the Java language level, and impact in the pthreads library or **_check_lock** is locking at the system level, which might or might not be associated with Java locks. For example, libpthreads.a and **_check_lock** are active for lock contention that is related to malloc on AIX.

Example: B-13 AIX profile excerpt showing kernel and shared library ticks

Total Ticks For All Processes (KERNEL) = 690

Subroutine	Ticks	%	Source	Address	Bytes
=====	=====	=====	=====	=====	=====
._check_lock	240	5.71	low.s	3420	40

Shared Object	Ticks	%	Address	Bytes
=====	=====	=====	=====	=====
libj9jit24.so	1157	27.51	900000003e81240	5c8878
libj9gc24.so	510	12.13	900000004534200	91d66
/usr/lib/libpthreads.a[shr_xpg5_64.o]	175	4.16	900000000b83200	30aa0

Profile: libj9jit24.so

Total Ticks For All Processes (libj9jit24.so) = 1157

Subroutine	Ticks	%	Source	Address	Bytes
=====	=====	=====	=====	=====	=====
._jitMonitorEntry	1121	26.66	nathehp.s	549fc0	cc0

- **GC impact:** The impact of initializing new objects and of GC is shown in Example B-13 as the 12.13% of ticks in the libj9gc24.so shared object. This high GC impact is related to the excessive creation of Double objects in the sample program.
- **Java method execution:** In Example B-14, the profile shows the time that is spent in the ProfileTest class, which is broken down by method. Some methods appear more than one time in the breakdown because they are compiled multiple times at increasing optimization levels by the JIT compiler. Most of the ticks appear in the final highly optimized version of the **doWork()D** method, into which the **initialize()V** and **calculate()V** methods are inlined by the JIT compiler.

Example: B-14 AIX profile excerpt showing Java classes and methods

Total Ticks For All Processes (JAVA) = 1450

Class	Ticks	%
=====	=====	=====
ProfileTest	1401	33.32
java/util/Random	38	0.90
java/lang/Float	5	0.12
java/lang/Double	3	0.07
java/lang/Math	3	0.07

Profile: ProfileTest

Total Ticks For All Processes (ProfileTest) = 1401

Method	Ticks	%	Source	Address	Bytes
=====	=====	=====	=====	=====	=====
doWork()D	1385	32.94	ProfileTest.java	1107283bc	b54
doWork()D	6	0.14	ProfileTest.java	110725148	464
doWork()D	4	0.10	ProfileTest.java	110726e3c	156c
initialize()V	3	0.07	ProfileTest.java	1107262dc	b4c
calculate()V	2	0.05	ProfileTest.java	110724400	144
initialize()V d04	1	0.02	ProfileTest.java	1107255c4	

Example B-15 contains a shell program to collect a profile on Linux using **OProfile**. The resulting profile might be similar to the previous example profile on AIX, indicating substantial time in spin locking and in GC. Depending on some specifics of the Linux system, however, the locking impact can appear in routines in the `libj9thr24.so` shared object, as compared to the AIX spin locking seen in `libj9jit24.so`.

You can also use the newer **perf** command to collect a profile on Linux. For example, to collect a profile for a running Java program you could use a command such as:

```
perf --events PM_RUN_CYC:500000 --separate-thread --separate-cpu --pid <java_pid> &
```

where `<java_pid>` was the process id of the running Java program. After a suitable length of profiling time you would enter **Ctrl-C** to terminate the profiling data collection. The **opreport** command would then be used to generate the profile report, in the same way as shown in Example B-15. As well as being simpler than the **OProfile** commands in Example B-15, one advantage of using **perf** in this way is that a profile of a single process does not require that the root userid be used. See the **OProfile** manual for more information.

In some cases, an environment variable setting might be necessary to indicate the location of the JVMTI library that is needed for running **OProfile** or **operf** with Java:

Linux 32-bit: **LD_LIBRARY_PATH=/usr/lib/oprofile**

Linux 64-bit: **LD_LIBRARY_PATH=/usr/lib64/oprofile**

Alternatively, you can specify the full path to the JVMTI library on the Java command line, such as this:

```
java -agentpath:/usr/lib/oprofile/libjvmti_oprofile.so
```

Example: B-15 Linux shell to collect a profile using OProfile

```
#!/bin/ksh

# Use --no-vmlinux if we either have a compressed kernel or do not care about the kernel symbols.
# Otherwise, use "opcontrol --vmlinux=/boot/vmlinux", for example.
opcontrol --no-vmlinux

# Stop data collection and remove daemon. Make sure we start from scratch.
opcontrol --shutdown

# Load the Oprofile module if required and makes the Oprofile driver interface available.
opcontrol --init

# Clear out data from current session.
# opcontrol --reset

# Select the performance counter that counts non-idle cycles and generates a sample after 500,000
# such events.
opcontrol -e PM_RUN_CYC_GRP1:500000

# Start the daemon for data collection.
opcontrol --start

# Run the Java program. "-agentlib:jvmti_oprofile" allows Oprofile to resolve the jitted methods.
java -Xms10m -Xmx10m -agentlib:jvmti_oprofile ProfileTest

# Stop data collection.
opcontrol --stop

# Flush the collected profiling data.
opcontrol --dump

# Generate a summary report at the module level.
opreport > ProfileTest_summary.log

# Generate a long report at the function level.
opreport -l > ProfileTest_long.log
```

Locking analysis

Locking bottlenecks are fairly common in Java applications. Collect locking information to identify any bottlenecks, and then take appropriate steps to eliminate the problems. A common case is when older java/util classes, such as Hashtable, do not scale well and cause a locking bottleneck. An easy solution is to use java/util/concurrent classes instead, such as ConcurrentHashMap.

Locking can be at the Java code level or at the system level. Java Lock Monitor is an easy to use tool that identifies locking bottlenecks at the Java language level or in internal JVM locking. A profile that is slowing a significant fraction of time in kernel locking routines indicates that system level locking that might be related to an underlying Java locking issue. Other AIX tools, such as **sp1at**, are helpful in diagnosing locking problems at the system level.

Always evaluate locking in the largest required scalability configuration (the largest number of cores).

Java Lock Monitor

The Java Lock Monitor is a valuable tool to deal with concurrency and synchronization in multi-threaded applications. The JLM can provide detailed information, such as how contested every monitor in the application is, how often a particular thread acquires a particular monitor, and how often a monitor is reacquired by a thread that already owns it. The locks that are surveyed by the JLM include both application locks and locks used internally by the JVM, such as GC locks. These statistics can be used to make decisions about GC policies, lock reservation, and so on, to make optimal usage of processing resources. For more information about the Java Lock Monitor, see *Java diagnostics, IBM style, Part 3: Diagnosing synchronization and locking problems with the Lock Analyzer for Java*, available here:

<http://www.ibm.com/developerworks/library/j-ibmtools3/>

Also, see “Hot method or routine analysis” on page 223.

Thread state analysis

Multi-threaded Java applications, especially applications that are running on top of WebSphere Application Server, often have many threads that might be blocked or waiting on locks, database operations, or file system operations. A powerful analysis technique is to look at the state of the threads to diagnose performance issues.

Always evaluate thread state analysis in the largest required scalability configuration (the largest number of cores).

IBM Whole-system Analysis of Idle Time (WAIT)

IBM Whole-system Analysis of Idle Time (WAIT) is a lightweight tool to assess various performance issues that range from GC to lock contention to file system bottlenecks and database bottlenecks, to client delays and authentication server delays, and more, including traditional performance issues, such as identifying hot methods.

WAIT was originally developed for Java and Java Platform, Enterprise Edition workloads, but a beta version that works with C/C++ native code is also available. The WAIT diagnostic capabilities are not limited to traditional Java bottlenecks such as GC problems or hot methods. WAIT employs an expert rule system to look at how Java code communicates with the wider world to provide a high-level view of system and application bottlenecks.

WAIT is also agentless (relying on `javacores`, `ps`, `vmstat`, and similar information, all of which are subject to availability). For example, WAIT produces a report with whatever subset of data can be extracted on a machine. Getting `javacores`, `ps`, and `vmstat` data almost never requires a change to command lines, environment variables, and so on.

Output is viewed in a browser such as Firefox, Chrome, Safari, and Internet Explorer, and assuming one has a browser, no additional installation is needed to view the WAIT output. Reports are interactive, and clicking different elements reveals more information. Manuals, animated demonstrations, and sample reports are also available on the WAIT website.

For more information about WAIT, go to this website:

<http://wait.researchlabs.ibm.com>

This site also has sample input files for WAIT, so users can try out the data analysis and visualization aspects without collecting any data.



Performance Optimization and Tuning Techniques for IBM Processors, including IBM POWER8

(0.5" spine)

0.475" <-> 0.873"

250 <-> 459 pages



Performance Optimization and Tuning Techniques for IBM Processors, including IBM POWER8



Redbooks®

Learn optimization strategies for the new IBM POWER8 processor

Apply strategies to IBM POWER7 and IBM POWER6 processors

Optimize code performance in POWER environments

This IBM Redbooks publication focuses on gathering the correct technical information, and laying out simple guidance for optimizing code performance on IBM POWER8 systems that run the AIX, IBM i, or Linux operating systems. There is much straightforward performance optimization that can be performed with a minimum of effort and without extensive previous experience or in-depth knowledge.

The POWER8 processor contains many new and important performance features, such as support for eight hardware threads in each core and support for transactional memory. POWER8 is a strict superset of IBM POWER7+, and so all of the performance features of POWER7+, such as multiple page sizes, also appear in POWER8. Much of the technical information and guidance for optimizing performance on POWER8 presented in this guide also applies to POWER7+ and earlier processors, except where the guide explicitly indicates that a feature is new in POWER8.

This guide strives to focus on optimizations that tend to be positive across a broad set of IBM POWER processor chips and systems. Specific guidance is given for the POWER8 processor; however, the general guidance is applicable to the IBM POWER7+, IBM POWER7, IBM POWER6, IBM POWER5, and even to earlier processors.

This guide is directed to personnel who are responsible for performing migration and implementation activities on IBM POWER8-based servers. This includes system administrators, system architects, network administrators, information architects, and database administrators (DBAs).

**INTERNATIONAL
TECHNICAL
SUPPORT
ORGANIZATION**

**BUILDING TECHNICAL
INFORMATION BASED ON
PRACTICAL EXPERIENCE**

IBM Redbooks are developed by the IBM International Technical Support Organization. Experts from IBM, Customers and Partners from around the world create timely technical information based on realistic scenarios. Specific recommendations are provided to help you implement IT solutions more effectively in your environment.

**For more information:
ibm.com/redbooks**

SG24-8171-00

ISBN 073843972X