Solaris Modular Debugger Guide

Alta Elstad

Solaris Modular Debugger Guide

by Alta Elstad

The contents of this Documentation are subject to the Public Documentation License Version 1.01 (the "License"); you may only use this Documentation if you comply with the terms of this License. A copy of the License is available at http://www.opensolaris.org/os/community/documentation/license.

Published June 2007 Copyright © 2007 Sun Microsystems, Inc.

Abstract

This book describes the Solaris Modular Debugger (MDB), which is a general purpose debugging tool for the Solaris operating system. The primary feature of MDB is its extensibility. This book describes how to use MDB to debug complex software systems, with a particular emphasis on the facilities available for debugging the Solaris kernel and associated device drivers and modules. The book also includes a complete reference for and discussion of the MDB language syntax, debugger features, and MDB Module Programming API.

Sun Microsystems, Inc. has intellectual property rights relating to technology embodied in the product that is described in this document. In particular, and without limitation, these intellectual property rights may include one or more U.S. patents or pending patent applications in the U.S. and in other countries.

U.S. Government Rights – Commercial software. Government users are subject to the Sun Microsystems, Inc. standard license agreement and applicable provisions of the FAR and its supplements.

This distribution may include materials developed by third parties.

Parts of the product may be derived from Berkeley BSD systems, licensed from the University of California. UNIX is a registered trademark in the U.S. and other countries, exclusively licensed through X/Open Company, Ltd.

Sun, Sun Microsystems, the Sun logo, the Solaris logo, the Java Coffee Cup logo, docs.sun.com, Java, and Solaris are trademarks or registered trademarks of Sun Microsystems, Inc. in the U.S. and other countries. All SPARC trademarks are used under license and are trademarks or registered trademarks of SPARC International, Inc. in the U.S. and other countries. Products bearing SPARC trademarks are based upon an architecture developed by Sun Microsystems, Inc.

The OPEN LOOK and SunTM Graphical User Interface was developed by Sun Microsystems, Inc. for its users and licensees. Sun acknowledges the pioneering efforts of Xerox in researching and developing the concept of visual or graphical user interfaces for the computer industry. Sun holds a non-exclusive license from Xerox to the Xerox Graphical User Interface, which license also covers Sun's licensees who implement OPEN LOOK GUIs and otherwise comply with Sun's written license agreements.

Products covered by and information contained in this publication are controlled by U.S. Export Control laws and may be subject to the export or import laws in other countries. Nuclear, missile, chemical or biological weapons or nuclear maritime end uses or end users, whether direct or indirect, are strictly prohibited. Export or reexport to countries subject to U.S. embargo or to entities identified on U.S. export exclusion lists, including, but not limited to, the denied persons and specially designated nationals lists is strictly prohibited.

DOCUMENTATION IS PROVIDED "AS IS" AND ALL EXPRESS OR IMPLIED CONDITIONS, REPRESENTATIONS AND WARRANTIES, INCLUDING ANY IMPLIED WARRANTY OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE OR NON-INFRINGEMENT, ARE DISCLAIMED, EXCEPT TO THE EXTENT THAT SUCH DISCLAIMERS ARE HELD TO BE LEGALLY INVALID.

Sun Microsystems, Inc. détient les droits de propriété intellectuelle relatifs à la technologie incorporée dans le produit qui est décrit dans ce document. En particulier, et ce sans limitation, ces droits de propriété intellectuelle peuvent inclure un ou plusieurs brevets américains ou des applications de brevet en attente aux Etats-Unis et dans d'autres pays.

Cette distribution peut comprendre des composants développés par des tierces personnes.

Certaines composants de ce produit peuvent être dérivées du logiciel Berkeley BSD, licenciés par l'Université de Californie. UNIX est une marque déposée aux Etats-Unis et dans d'autres pays; elle est licenciée exclusivement par X/Open Company, Ltd.

Sun, Sun Microsystems, le logo Sun, le logo Solaris, le logo Java Coffee Cup, docs.sun.com, Java et Solaris sont des marques de fabrique ou des marques déposées de Sun Microsystems, Inc. aux Etats-Unis et dans d'autres pays. Toutes les marques SPARC sont utilisées sous licence et sont des marques de fabrique ou des marques déposées de SPARC International, Inc. aux Etats-Unis et dans d'autres pays. Les produits portant les marques SPARC sont basés sur une architecture développée par Sun Microsystems, Inc.

L'interface d'utilisation graphique OPEN LOOK et Sun a été développée par Sun Microsystems, Inc. pour ses utilisateurs et licenciés. Sun reconnaît les efforts de pionniers de Xerox pour la recherche et le développement du concept des interfaces d'utilisation visuelle ou graphique pour l'industrie

de l'informatique. Sun détient une licence non exclusive de Xerox sur l'interface d'utilisation graphique Xerox, cette licence couvrant également les licenciés de Sun qui mettent en place l'interface d'utilisation graphique OPEN LOOK et qui, en outre, se conforment aux licences écrites de Sun.

Les produits qui font l'objet de cette publication et les informations qu'il contient sont régis par la legislation américaine en matière de contrôle des exportations et peuvent être soumis au droit d'autres pays dans le domaine des exportations et importations. Les utilisations finales, ou utilisateurs finaux, pour des armes nucléaires, des missiles, des armes chimiques ou biologiques ou pour le nucléaire maritime, directement ou indirectement, sont strictement interdites. Les exportations ou réexportations vers des pays sous embargo des Etats-Unis, ou vers des entités figurant sur les listes d'exclusion d'exportation américaines, y compris, mais de manière non exclusive, la liste de personnes qui font objet d'un ordre de ne pas participer, d'une façon directe ou indirecte, aux exportations des produits ou des services qui sont régis par la legislation américaine en matière de contrôle des exportations et la liste de ressortissants spécifiquement designés, sont rigoureusement interdites.

LA DOCUMENTATION EST FOURNIE "EN L'ETAT" ET TOUTES AUTRES CONDITIONS, DECLARATIONS ET GARANTIES EXPRESSES OU TACITES SONT FORMELLEMENT EXCLUES, DANS LA MESURE AUTORISEE PAR LA LOI APPLICABLE, Y COMPRIS NOTAMMENT TOUTE GARANTIE IMPLICITE RELATIVE A LA QUALITE MARCHANDE, A L'APTITUDE A UNE UTILISATION PARTICULIERE OU A L'ABSENCE DE CONTREFACON.

Table of Contents

Preface	
Who Should Use This Book	X
How This Book Is Organized	X
Related Books and Papers	. xi
Accessing Sun Documentation Online	xii
What Typographic Conventions Mean	xii
Shell Prompts in Command Examples	xii
1. Modular Debugger Overview	
Introduction	1
MDB Features	
Using MDB	2
Future Enhancements	
2. Debugger Concepts	
Building Blocks	
Modularity	
3. Language Syntax	
Syntax	
Commands	
Comments	
Arithmetic Expansion	
Unary Operators	
Binary Operators	
Quoting	
Shell Escapes	
1	
Variables	
Symbol Name Resolution	
dcmd and Walker Name Resolution	
dcmd Pipelines	
Formatting demds	
4. Interaction	
Command Re-entry	
In-line Editing	
Keyboard Shortcuts	
Output Pager	
Signal Handling	
5. Built-in Commands	
Built-in dcmds	
6. Execution Control	
Execution Control	
Event Callbacks	
Thread Support	30
Built-in dcmds	30
Interaction With exec	. 35
Interaction with Job Control	. 35
Process Attach and Release	. 36
7. Kernel Execution Control	. 37
Booting, Loading, and Unloading	
Terminal Handling	
Debugger Entry	
Processor-Specific Features	
8. Kernel Debugging Modules	

Generic Kernel Debugging Support (genunix)	
Kernel Memory Allocator	
File Systems	41
Virtual Memory	. 42
CPUs and the Dispatcher	43
Device Drivers and DDI Framework	. 43
STREAMS	45
Networking	. 46
Files, Processes, and Threads	
Synchronization Primitives	
Cyclics	
Task Queues	
Error Queues	
Configuration	
Interprocess Communication Debugging Support (ipc)	
dcmds	
Walkers	
Loopback File System Debugging Support (lofs)	
dcmds	
Walkers	
Internet Protocol Module Debugging Support (ip)	
dcmds	
Walkers	
Kernel Runtime Link Editor Debugging Support (krtld)	
demds	
Walkers	
USB Framework Debugging Support (uhci)	
dcmds	
Walkers	
USB Framework Debugging Support (usba)	
dcmds	
Walkers	
x86 Platform Debugging Support (unix)	
dcmds	
Walkers	
sun4u Platform Debugging Support (unix)	
dcmds	
Walkers	. 55
9. Debugging With the Kernel Memory Allocator	
Getting Started: Creating a Sample Crash Dump	56
Setting kmem_flags	56
Forcing a Crash Dump	56
Starting MDB	57
Allocator Basics	. 57
Buffer States	. 58
Transactions	. 58
Sleeping and Non-Sleeping Allocations	58
Kernel Memory Caches	58
Kernel Memory Caches	
Detecting Memory Corruption	
Freed Buffer Checking: 0xdeadbeef	
Redzone: Oxfeedface	
Uninitialized Data: 0xbaddcafe	
Associating Panic Messages With Failures	

Memory Allocation Logging	
Buftag Data Integrity	. 65
The bufctl Pointer	. 65
Advanced Memory Analysis	. 66
Finding Memory Leaks	. 67
Finding References to Data	. 67
Finding Corrupt Buffers With ::kmem_verify	68
Allocator Logging Facility	
10. Module Programming API	
Debugger Module Linkage	
_mdb_init	
_mdb_fini	
Dcmd Definitions	
Walker Definitions	
API Functions	
mdb_pwalk	
mdb_walk	
mdb pwalk dcmd	
 -	
mdb_walk_dcmd	
mdb_call_dcmd	
mdb_layered_walk	
mdb_add_walker	
mdb_remove_walker	
mdb_vread and mdb_vwrite	
mdb_fread and mdb_fwrite	
mdb_pread and mdb_pwrite	. 78
mdb_readstr	. 78
mdb_writestr	. 79
mdb_readsym	. 79
mdb_writesym	. 79
mdb_readvar and mdb_writevar	. 79
mdb_lookup_by_name and mdb_lookup_by_obj	
mdb_lookup_by_addr	
mdb getopts	
mdb_strtoull	
mdb alloc, mdb zalloc and mdb free	
mdb_printf	
mdb snprintf	
mdb_warn	
mdb flush	
mdb_nhconvert	
mdb_dumpptr and mdb_dump64	
mdb_one_bit	
mdb_inval_bits	
mdb_inc_indent and mdb_dec_indent	
mdb_eval	
mdb_set_dot and mdb_get_dot	
mdb_get_pipe	
mdb_set_pipe	
mdb_get_xdata	
Additional Functions	
A. Options	
Summary of Command-line Options	
Operands	. 94

Solaris Modular Debugger Guide

Exit Status	. 94
Environment Variables	. 94
B. Notes	. 95
Warnings	. 95
Use of the Error Recovery Mechanism	. 95
Use of the Debugger to Modify the Live Operating System	. 95
Use of kmdb to Stop the Live Operating System	. 95
Notes	. 95
Limitations on Examining Process Core Files	. 95
Limitations on Examining Crash Dump Files	. 95
Relationship Between 32-bit and 64-bit Debugger	. 96
Limitations on Memory Available to kmdb	. 96
Developer Information	. 96
C. Transition From adb and kadb	. 97
Command-line Options	. 97
Syntax	97
Watchpoint Length Specifier	. 97
Address Map Modifier	. 98
Output	98
Deferred Breakpoints	. 98
I/O Port Access	. 98
D. Transition From crash	. 99
Command-line Options	. 99
Input in MDB	. 99
Functions	99

List of Figures

2.1. MDB architecture	. 4
9.1. The Redzone	
9.2. Sample kmem_alloc(9F) Buffer	63
9.3. Redzone Byte	63
9.4. Redzone Byte at the Beginning of the Redzone	64
9.5. Extra Debugging Data in the Buftag	65
10.1. Sample Walker	75

List of Tables

1. Typographic Conventions	xiii
2. Shell Prompts	xiv
D.1. Radix Specifiers	99

Preface

The Modular Debugger (MDB) is a highly extensible, general purpose debugging tool for the SolarisTM Operating System. The *Solaris Modular Debugger Guide* describes how to use MDB to debug complex software systems, with a particular emphasis on the facilities available for debugging the Solaris kernel and associated device drivers and modules. It also includes a complete reference for and discussion of the MDB language syntax, debugger features, and MDB Module Programming API.

Note

This Solaris release supports systems that use the SPARC® and x86 families of processor architectures: UltraSPARC®, SPARC64, AMD64, Pentium, and Xeon EM64T. The supported systems appear in the *Solaris OS Hardware Compatibility Lists* at http://www.sun.com/bigadmin/hcl. This document cites any implementation differences between the platform types.

In this document the term "x86" refers to 64-bit and 32-bit systems manufactured using processors compatible with the AMD64 or Intel Xeon/Pentium product families. For supported systems, see the *Solaris OS Hardware Compatibility Lists*.

Who Should Use This Book

If you were a detective and were investigating at the scene of a crime, you might interview the witnesses and ask them to describe what happened and who they saw. However, if there were no witnesses or these descriptions proved insufficient, you might consider collecting fingerprints and forensic evidence that could be examined for DNA to help solve the case. Often, software program failures divide into analogous categories: problems that can be solved with source-level debugging tools, and problems that require low-level debugging facilities, examination of core files, and knowledge of assembly language to diagnose and correct. MDB is a debugger designed to facilitate analysis of this second class of problems.

It might not be necessary to use MDB in every case, just as a detective doesn't need a microscope and DNA evidence to solve every crime. However, when programming a complex low-level software system such as an operating system, these situations can occur frequently. As a result, MDB is designed as a debugging framework that allows you to construct your own custom analysis tools to aid in the diagnosis of these problems. MDB also provides a powerful set of built-in commands that allow you to analyze the state of your program at the assembly language level.

If you are not familiar with assembly language programming and debugging, Related Books and Papers provides references to materials that you might find useful.

You should also disassemble various functions of interest in the programs you will be debugging in order to familiarize yourself with the relationship between your program's source code and the corresponding assembly language code. If you are planning to use MDB for debugging Solaris kernel software, you should read carefully Chapter 8, Kernel Debugging Modules and Chapter 9, Debugging With the Kernel Memory Allocator. These chapters provide more detailed information on the MDB commands and facilities provided for debugging Solaris kernel software.

How This Book Is Organized

Chapter 1, Modular Debugger Overview provides an overview of the debugger. This chapter is intended for all users.

Chapter 2, Debugger Concepts describes the MDB architecture and explains the terminology for the debugger concepts used throughout this book. This chapter is intended for all users.

Chapter 3, Language Syntax describes the syntax, operators and evaluation rules for the MDB language. This chapter is intended for all users.

Chapter 4, Interaction describes the MDB interactive command-line editing facilities and output pager. This chapter is intended for all users.

Chapter 5, Built-in Commands describes the set of built-in debugger commands that are always available. This chapter is intended for all users.

Chapter 6, Execution Control describes the MDB facilities for controlling the execution of live running programs. This chapter is intended for application developers and device driver developers. Execution control features might also be useful for system administrators.

Chapter 7, Kernel Execution Control describes the MDB facilities for controlling the execution of the live operating system kernel that are specific to **kmdb**. This chapter is intended for operating system kernel developers and device driver developers.

Chapter 8, Kernel Debugging Modules describes the set of loadable debugger commands that are provided for debugging the Solaris kernel. This chapter is intended for users who intend to examine Solaris kernel crash dumps and for kernel software developers.

Chapter 9, Debugging With the Kernel Memory Allocator describes the debugging features of the Solaris kernel memory allocator and the MDB commands provided to take advantage of these features. This chapter is intended for advanced programmers and kernel software developers.

Chapter 10, Module Programming API describes the facilities for writing loadable debugger modules. This chapter is intended for advanced programmers and software developers who intend to develop custom debugging support for MDB.

Appendix A, Options provides a reference for MDB command-line options.

Appendix B, Notes provides warnings and notes about using the debugger.

Appendix C, Transition From adb and kadb provides a reference for **adb** commands and their MDB equivalents. The **adb** command is implemented by **mdb**.

Appendix D, Transition From crash provides a reference for **crash** commands and their MDB equivalents. The **crash** command is no longer present in Solaris.

Related Books and Papers

These books and papers are recommended and related to the tasks that you need to perform:

- Vahalia, Uresh. UNIX Internals: The New Frontiers. Prentice Hall, 1996. ISBN 0-13-101908-2
- Mauro, Jim and McDougall, Richard. *Solaris Internals: Core Kernel Components*. Sun Microsystems Press, 2001. ISBN 0-13-022496-0
- The SPARC Architecture Manual, Version 9. Prentice Hall, 1998. ISBN 0-13-099227-5
- The SPARC Architecture Manual, Version 8. Prentice Hall, 1994. ISBN 0-13-825001-4

- *Pentium Pro Family Developer's Manual, Volumes 1-3.* Intel Corporation, 1996. ISBN 1-55512-259-0 (Volume 1), ISBN 1-55512-260-4 (Volume 2), ISBN 1-55512-261-2 (Volume 3)
- Bonwick, Jeff. *The Slab Allocator: An Object-Caching Kernel Memory Allocator*. Proceedings of the Summer 1994 Usenix Conference, 1994. ISBN 9–99–452010–5
- SPARC Assembly Language Reference Manual. Sun Microsystems, 1998.
- x86 Assembly Language Reference Manual. Sun Microsystems, 1998.
- Writing Device Drivers. Sun Microsystems, 2000.
- STREAMS Programming Guide. Sun Microsystems, 2000.
- Solaris 64-bit Developer's Guide. Sun Microsystems, 2000.
- Linker and Libraries Guide. Sun Microsystems, 2000.

Accessing Sun Documentation Online

The docs.sun.comSM Web site enables you to access Sun technical documentation online. You can browse the docs.sun.com archive or search for a specific book title or subject. The URL is http://docs.sun.com.

What Typographic Conventions Mean

The following table describes the typographic changes used in this book.

Table 1. Typographic Conventions

Typeface or Symbol	Meaning	Example
AaBbCc123	The names of commands, files, and directories; on-screen computer output	Edit your .login file. Use ls a to list all files. machine_name% you have mail.
AaBbCc123	What you type, contrasted with on-screen computer output	machine_name% su Password:
AaBbCc123	Command-line placeholder: replace with a real name or value	To delete a file, type rm filename.
AaBbCc123	Book titles, new words, or terms, or words to be emphasized.	Read Chapter 6 in <i>User's Guide</i> . These are called <i>class</i> options. You must be <i>root</i> to do this.

Shell Prompts in Command Examples

The following table shows the default system prompt and superuser prompt for the C shell, Bourne shell, and Korn shell, as well as the MDB debugger prompt.

Table 2. Shell Prompts

Shell	Prompt
C shell prompt	machine_name%
C shell superuser prompt	machine_name#
Bourne shell and Korn shell prompt	\$
Bourne shell and Korn shell superuser prompt	#
mdb prompt	>
kmdb prompt	[cpu]>

Chapter 1. Modular Debugger Overview

The Modular Debugger (MDB) is a general purpose debugging tool for Solaris whose primary feature is its extensibility. This book describes how to use MDB to debug complex software systems, with a particular emphasis on the facilities available for debugging the Solaris kernel and associated device drivers and modules. The book also includes a complete reference for and discussion of the MDB language syntax, debugger features, and MDB Module Programming API.

Introduction

Debugging is the process of analyzing the execution and state of a software program in order to remove defects. Traditional debugging tools provide facilities for execution control so that programmers can reexecute programs in a controlled environment and display the current state of program data or evaluate expressions in the source language used to develop the program. Unfortunately, these techniques are often inappropriate for debugging complex software systems such as:

- · An operating system, where bugs might not be reproducible and program state is massive and distributed
- Programs that are highly optimized or have had their debug information removed
- · Programs that are themselves low-level debugging tools
- Customer situations where the developer can only access post-mortem information

MDB is a tool that provides a completely customizable environment for debugging these programs and scenarios, including a dynamic module facility that programmers can use to implement their own debugging commands to perform program-specific analysis. Each MDB module can be used to examine the program in several different contexts, including live and post-mortem. The Solaris Operating System includes a set of MDB modules designed to aid programmers in debugging the Solaris kernel and related device drivers and kernel modules. Third-party developers might find it useful to develop and deliver their own debugging modules for supervisor or user software.

MDB Features

MDB provides an extensive collection of features for analyzing the Solaris kernel and other target programs. You can:

- Perform post-mortem analysis of Solaris kernel crash dumps and user process core dumps: MDB
 includes a collection of debugger modules that facilitate sophisticated analysis of kernel and process
 state, in addition to standard data display and formatting capabilities. The debugger modules allow you
 to formulate complex queries to:
 - Locate all the memory allocated by a particular thread
 - Print a visual picture of a kernel STREAM
 - Determine what type of structure a particular address refers to
 - · Locate leaked memory blocks in the kernel
 - · Analyze memory to locate stack traces
- Use a first-class programming API to implement your own debugger commands and analysis tools
 without having to recompile or modify the debugger itself: In MDB, debugging support is implemented

as a set of loadable modules (shared libraries that the debugger can dlopen(3C)), each of which provides a set of commands that extends the capabilities of the debugger itself. The debugger in turn provides an API of core services, such as the ability to read and write memory and access symbol table information. MDB provides a framework for developers to implement debugging support for their own drivers and modules; these modules can then be made available for everyone to use.

- Learn to use MDB if you are already familiar with the legacy debugging tools **adb** and **crash**: MDB provides backward compatibility with these existing debugging solutions. The MDB language itself is designed as a superset of the adb language; all existing adb macros and commands work within MDB so developers who use adb can immediately use MDB without knowing any MDB-specific commands. MDB also provides commands that surpass the functionality available from the crash utility.
- Benefit from enhanced usability features. MDB provides a host of usability features, including:
 - · Command-line editing
 - Command history
 - · Built-in output pager
 - · Syntax error checking and handling
 - Online help
 - · Interactive session logging

Using MDB

MDB is available on Solaris systems as two commands that share common features: **mdb** and **kmdb**. You can use the **mdb** command interactively or in scripts to debug live user processes, user process core files, kernel crash dumps, the live operating system, object files, and other files. You can use the **kmdb** command to debug the live operating system kernel and device drivers when you also need to control and halt the execution of the kernel. To start **mdb**, execute the mdb(1) command. To start **kmdb**, boot the system as described in the kmdb(1) man page, or execute the **mdb** command with the K option.

Future Enhancements

MDB provides a stable foundation for developing advanced post-mortem analysis tools. Each Solaris operating system release includes additional MDB modules that provide even more sophisticated functionality for debugging the kernel and other software programs. You can use MDB to debug existing software programs, and develop your own modules to improve your ability to debug your own Solaris drivers and applications.

Chapter 2. Debugger Concepts

This section discusses the significant aspects of MDB's design, and the benefits derived from this architecture.

Building Blocks

The **target** is the program being inspected by the debugger. MDB currently provides support for the following types of targets:

- User processes
- User process core files
- Live operating system without kernel execution control (through /dev/kmem and /dev/ksyms)
- Live operating system with kernel execution control (through kmdb(1))
- · Operating system crash dumps
- User process images recorded inside an operating system crash dump
- ELF object files
- · Raw data files

Each target exports a standard set of properties, including one or more address spaces, one or more symbol tables, a set of load objects, and a set of threads. Figure 2–1 shows an overview of the MDB architecture, including two of the built-in targets and a pair of sample modules.

A debugger command, or **dcmd** (pronounced *dee-command*) in MDB terminology, is a routine in the debugger that can access any of the properties of the current target. MDB parses commands from standard input, then executes the corresponding dcmds. Each dcmd can also accept a list of string or numerical arguments, as shown in Syntax. MDB contains a set of built-in dcmds described in Chapter 5, Built-in Commands, that are always available. The programmer can also extend the capabilities of MDB itself by writing dcmds using a programming API provided with MDB.

A walker is a set of routines that describe how to walk, or iterate, through the elements of a particular program data structure. A walker encapsulates the data structure's implementation from dcmds and from MDB itself. You can use walkers interactively, or use them as a primitive to build other dcmds or walkers. As with dcmds, the programmer can extend MDB by implementing additional walkers as part of a debugger module.

A debugger module, or **dmod** (pronounced *dee-mod*), is a dynamically loaded library containing a set of dcmds and walkers. During initialization, MDB attempts to load dmods corresponding to the load objects present in the target. You can subsequently load or unload dmods at any time while running MDB. MDB provides a set of standard dmods for debugging the Solaris kernel.

A macro file is a text file containing a set of commands to execute. Macro files are typically used to automate the process of displaying a simple data structure. MDB provides complete backward compatibility for the execution of macro files written for adb. The set of macro files provided with the Solaris installation can therefore be used with either tool.

Figure 2.1. MDB architecture

This graphic describes the MDB's components: the MDB language and the MDB module API overlying the debugger engine.

Modularity

The benefit of MDB's modular architecture extends beyond the ability to load a module containing additional debugger commands. The MDB architecture defines clear interface boundaries between each of the layers shown in Figure 2–1. Macro files execute commands written in the MDB or adb language. Dcmds and walkers in debugger modules are written using the MDB Module API, and this forms the basis of an application binary interface that allows the debugger and its modules to evolve independently.

The MDB name space of walkers and dcmds also defines a second set of layers between debugging code that maximizes code sharing and limits the amount of code that must be modified as the target program itself evolves. For example, one of the primary data structures in the Solaris kernel is the list of proc_t structures representing active processes in the system. The ::ps dcmd must iterate over this list in order to produce its output. However, the code to iterate over the list is not in the ::ps dcmd, it is encapsulated in the genunix module's proc walker.

MDB provides both ::ps and ::ptree dcmds, but neither has any knowledge of how proc_t structures are accessed in the kernel. Instead, they invoke the proc walker programmatically and format the set of returned structures appropriately. If the data structure used for proc_t structures ever changed, MDB could provide a new proc walker and none of the dependent dcmds would need to change. The proc walker can also be accessed interactively using the ::walk dcmd in order to create novel commands as you work during a debugging session.

In addition to facilitating layering and code sharing, the MDB Module API provides dcmds and walkers with a single stable interface for accessing various properties of the underlying target. The same API functions are used to access information from user process or kernel targets, simplifying the task of developing new debugging facilities.

In addition, a custom MDB module can be used to perform debugging tasks in a variety of contexts. For example, you might want to develop an MDB module for a user program you are developing. Once you have done so, you can use this module when MDB examines a live process executing your program, a core dump of your program, or even a kernel crash dump taken on a system where your program was executing.

The Module API provides facilities for accessing the following target properties:

Address module API provides facilities for reading and writing data from the target's virtual address Spacespace. Functions for reading and writing using physical addresses are also provided for kernel debugging modules.

Symbolihe module API provides access to the static and dynamic symbol tables of the target's primary Table executable file, its runtime link-editor, and a set of load objects (shared libraries in a user process or loadable modules in the Solaris kernel).

External module API provides a facility for retrieving a collection of named external data buffers Data associated with the target. For example, MDB provides programmatic access to the proc(4) structures associated with a user process or user core file target.

In addition, you can use built-in MDB dcmds to access information about target memory mappings, load objects, register values, and control the execution of user process targets.

Chapter 3. Language Syntax

This chapter describes the MDB language syntax, operators, and rules for command and symbol name resolution.

Syntax

The debugger processes commands from standard input. If standard input is a terminal, MDB provides terminal editing capabilities. MDB can also process commands from macro files and from dcmd pipelines, described below. The language syntax is designed around the concept of computing the value of an expression (typically a memory address in the target), and applying a dcmd to that address. The current address location is referred to as dot, and "." is used to reference its value.

A metacharacter is one of the following characters:

A blank is a TAB or a SPACE. A word is a sequence of characters separated by one or more non-quoted metacharacters. Some of the metacharacters function only as delimiters in certain contexts, as described below. An identifier is a sequence of letters, digits, underscores, periods, or back quotes beginning with a letter, underscore, or period. Identifiers are used as the names of symbols, variables, dcmds, and walkers. Commands are delimited by a NEWLINE or semicolon (;).

A dcmd is denoted by one of the following words or metacharacters:

```
/ \ ? = > $character ::identifier
```

dcmds named by metacharacters or prefixed by a single \$ or : are provided as built-in operators, and implement complete compatibility with the command set of the legacy adb(1) utility. After a dcmd has been parsed, the /, \backslash , ?, =, >, \$, and : characters are no longer recognized as metacharacters until the termination of the argument list.

A simple-command is a dcmd followed by a sequence of zero or more blank-separated words. The words are passed as arguments to the invoked dcmd, except as specified under Arithmetic Expansion and Quoting. Each dcmd returns an exit status that indicates it was either successful, failed, or was invoked with invalid arguments.

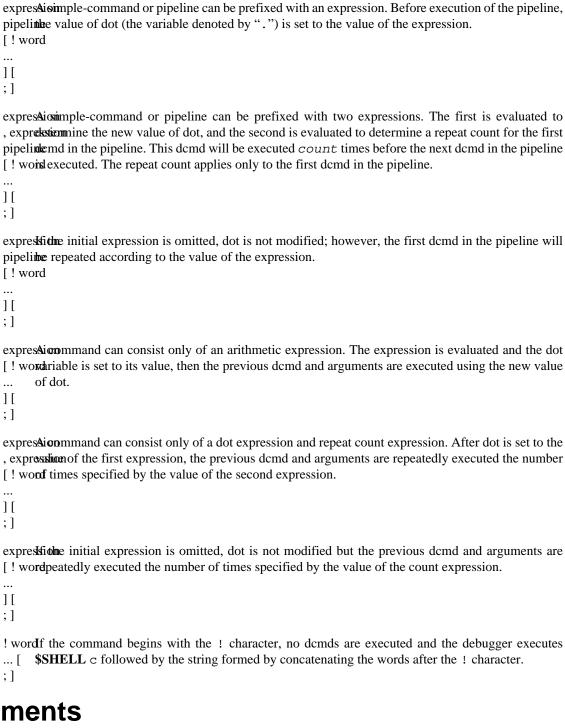
A pipeline is a sequence of one or more simple commands separated by |. Unlike the shell, dcmds in MDB pipelines are not executed as separate processes. After the pipeline has been parsed, each dcmd is invoked in order from left to right. Each dcmd's output is processed and stored as described in dcmd Pipelines. After the left-hand dcmd is complete, its processed output is used as input for the next dcmd in the pipeline. If any dcmd does not return a successful exit status, the pipeline is aborted.

An *expression* is a sequence of words that is evaluated to compute a 64-bit unsigned integer value. The words are evaluated using the rules described in Arithmetic Expansion.

Commands

A command is one of the following:

pipeline simple-command or pipeline can be optionally suffixed with the ! character, indicating that the [! wordebugger should open a pipe(2) and send the standard output of the last dcmd in the MDB pipeline to ... an external process created by executing \$SHELL c followed by the string formed by concatenating [the words after the ! character. For more details, refer to Shell Escapes. ;]



Comments

A word beginning with // causes that word and all the subsequent characters up to a NEWLINE to be ignored.

Arithmetic Expansion

Arithmetic expansion is performed when an MDB command is preceded by an optional expression representing a start address, or a start address and a repeat count. Arithmetic expansion can also be

performed to compute a numerical argument for a dcmd. An arithmetic expression can appear in an argument list enclosed in square brackets preceded by a dollar sign (\$[expression]), and will be replaced by the value of the expression.

Expressions can contain any of the following special words:

indicate binary values, 00 or 00 to indicate octal values, 0t or 0T to indicate decimal values, and 0x or 0X to indicate hexadecimal values (the default). The specified decimal floating point value, converted to its IEEE double-0[tT][0-9]+.[0-9]+ precision floating point representation 'ccccccc' The integer value computed by converting each character to a byte equal to its ASCII value. Up to eight characters can be specified in a character

constant. Characters are packed into the integer in reverse order (right-toleft), beginning at the least significant byte.

The specified integer value. Integer values can be prefixed with 0i or 0I to

<identifier The value of the variable named by identifier

identifier The value of the symbol named by identifier

(expression) The value of expression

The value of dot

& The most recent value of dot used to execute a dcmd

The value of dot incremented by the current increment

The value of dot decremented by the current increment

The increment is a global variable that stores the total bytes read by the last formatting dcmd. For more information on the increment, refer to the discussion of Formatting dcmds.

Unary Operators

integer

Unary operators are right associative and have higher precedence than binary operators. The unary operators are:

#expression Logical negation

~expression Bitwise complement

-expression Integer negation

%expression Value of a pointer-sized quantity at the object file location corresponding

to virtual address expression in the target's virtual address space

Value of a char-, short-, int-, or long-sized quantity at the object file %/[csil]/expression

location corresponding to virtual address expression in the target's

virtual address space

%/[1248]/expression Value of a one-, two-, four-, or eight-byte quantity at the object file

location corresponding to virtual address expression in the target's

virtual address space

*expression Value of a pointer-sized quantity at virtual address expression in the

target's virtual address space

*/[csil]/expression Value of a char-, short-, int-, or long-sized quantity at virtual address

expression in the target's virtual address space

*/[1248]/expression Value of a one-, two-, four-, or eight-byte quantity at virtual address

expression in the target's virtual address space

Binary Operators

Binary operators are left associative and have lower precedence than unary operators. The binary operators, in order of precedence from highest to lowest, are:

- * Integer multiplication
- % Integer division
- # Left-hand side rounded up to next multiple of right-hand side
- + Integer addition
- Integer subtraction
- << Bitwise shift left
- >> Bitwise shift right
- == Logical equality
- ! = Logical inequality
- & Bitwise AND
- A Bitwise exclusive OR
- Bitwise inclusive OR

Quoting

Each metacharacter described previously (see Syntax) terminates a word unless quoted. Characters can be quoted (forcing MDB to interpret each character as itself without any special significance) by enclosing them in a pair of single (') or double (") quotation marks. A single quote cannot appear within single quotes. Inside double quotes, MDB recognizes the C programming language character escape sequences.

Shell Escapes

The ! character can be used to create a pipeline between an MDB command and the user's shell. Shell escapes are only available when using **mdb** and not **kmdb**. If the \$SHELL environment variable is set, MDB will **fork** and **exec** this program for shell escapes; otherwise **/bin/sh** is used. The shell is invoked with the c option followed by a string formed by concatenating the words after the ! character.

The ! character takes precedence over all other metacharacters, except semicolon (;) and NEWLINE. After a shell escape is detected, the remaining characters up to the next semicolon or NEWLINE are passed "as is" to the shell. The output of shell commands cannot be piped to MDB dcmds. Commands executed by a shell escape have their output sent directly to the terminal, not to MDB.

Variables

A variable is a variable name, a corresponding integer value, and a set of attributes. A variable name is a sequence of letters, digits, underscores, or periods. A variable can be assigned a value using the > dcmd or ::typeset dcmd, and its attributes can be manipulated using the ::typeset dcmd. Each variable's value is represented as a 64-bit unsigned integer. A variable can have one or more of the following attributes: read-only (cannot be modified by the user), persistent (cannot be unset by the user), and tagged (user-defined indicator).

The following variables are defined as persistent:

0	Most recent value printed using the $/$, \setminus , ?, or = dcmd
9	Most recent count used with the \$< dcmd
b	Virtual address of the base of the data section
cpuid	The CPU identifier corresponding to the CPU on which kmdb is currently executing.
d	Size of the data section in bytes
e	Virtual address of the entry point
hits	The count of the number of times the matched software event specifier has been matched. See Event Callbacks.
m	Initial bytes (magic number) of the target's primary object file, or zero if no object file has been read yet
t	Size of the text section in bytes

t Size of the text section in bytes

thread The thread identifier of the current representative thread. The value of the identifier depends on the threading model used by the current target. See Thread Support.

In addition, the MDB kernel and process targets export the current values of the representative thread's register set as named variables. The names of these variables depend on the target's platform and instruction set architecture.

Symbol Name Resolution

As explained in Syntax, a symbol identifier present in an expression context evaluates to the value of this symbol. The value typically denotes the virtual address of the storage associated with the symbol in the target's virtual address space. A target can support multiple symbol tables including, but not limited to,

- Primary executable symbol table
- Primary dynamic symbol table
- Runtime link-editor symbol table
- Standard and dynamic symbol tables for each of a number of load objects (such as shared libraries in a user process, or kernel modules in the Solaris kernel)

The target typically searches the primary executable's symbol tables first, then one or more of the other symbol tables. Notice that ELF symbol tables contain only entries for external, global, and static symbols; automatic symbols do not appear in the symbol tables processed by MDB.

Additionally, MDB provides a private user-defined symbol table that is searched prior to any of the target symbol tables. The private symbol table is initially empty, and can be manipulated using the ::nmadd and ::nmdel dcmds.

The ::nm P option can be used to display the contents of the private symbol table. The private symbol table allows the user to create symbol definitions for program functions or data that were either missing from the original program or stripped out. These definitions are then used whenever MDB converts a symbolic name to an address, or an address to the nearest symbol.

Because targets contain multiple symbol tables, and each symbol table can include symbols from multiple object files, different symbols with the same name can exist. MDB uses the backquote " ` " character as a symbol-name scoping operator to allow the programmer to obtain the value of the desired symbol in this situation.

You can specify the scope used to resolve a symbol name as either: <code>object`name</code>, or <code>file`name</code>, or <code>object`file`name</code>. The object identifier refers to the name of a load object. The file identifier refers to the basename of a source file that has a symbol of type <code>STT_FILE</code> in the specified object's symbol table. The object identifier's interpretation depends on the target type.

The MDB kernel target expects *object* to specify the base name of a loaded kernel module. For example, the symbol name:

```
specfs`_init
```

evaluates to the value of the _init symbol in the specfs kernel module.

The **mdb** process target expects *object* to specify the name of the executable or of a loaded shared library. It can take any of the following forms:

- Exact match (that is, a full pathname): /usr/lib/libc.so.1
- Exact basename match: libc.so.1
- Initial basename match up to a ``." suffix: libc.so or libc
- Literal string a . out which is accepted as an alias for the executable

The process target will also accept any of the four forms described above preceded by an optional link-map id (lmid). The lmid prefix is specified by an initial LM followed by the link-map id in hexadecimal followed by an additional backquote. For example, the symbol name:

```
LM0`libc.so.1`_init
```

will evaluate to the value of the _init symbol in the libc.so.1 library that is loaded on link-map 0 (LM_ID_BASE). The link-map specifier may be necessary to resolve symbol naming conflicts in the event that the same library is loaded on more than one link map. For more information on link maps, refer to the Linker and Libraries Guide and the dlopen(3C) man page. Link-map identifiers will be displayed when symbols are printed according to the setting of the showlmid option, as described under Summary of Command-line Options.

In the case of a naming conflict between symbols and hexadecimal integer values, MDB attempts to evaluate an ambiguous token as a symbol first, before evaluating it as an integer value. For example, the token f can refer either to the decimal integer value 15 specified in hexadecimal (the default base), or to a global variable named f in the target's symbol table. If a symbol with an ambiguous name is present, the integer value can be specified by using an explicit 0x or 0x prefix.

dcmd and Walker Name Resolution

As described earlier, each MDB dmod provides a set of dcmds and walkers. dcmds and walkers are tracked in two distinct, global namespaces. MDB also keeps track of a dcmd and walker namespace associated with each dmod. Identically named dcmds or walkers within a given dmod are not allowed: a dmod with this type of naming conflict will fail to load.

Name conflicts between dcmds or walkers from different dmods are allowed in the global namespace. In the case of a conflict, the first dcmd or walker with that particular name to be loaded is given precedence in the global namespace. Alternate definitions are kept in a list in load order.

The backquote character "` " can be used in a dcmd or walker name as a scoping operator to select an alternate definition. For example, if dmods m1 and m2 each provide a dcmd d, and m1 is loaded prior to m2, then:

```
::d Executes m1's definition of d
```

::m1 d Executes m1's definition of d

:: m2 `d Executes m2 's definition of d

If module m1 were now unloaded, the next dcmd on the global definition list (m2 `d) would be promoted to global visibility. The current definition of a dcmd or walker can be determined using the ::which dcmd, described below. The global definition list can be displayed using the ::which v option.

dcmd Pipelines

dcmds can be composed into a pipeline using the | operator. The purpose of a pipeline is to pass a list of values, typically virtual addresses, from one dcmd or walker to another. Pipeline stages might be used to map a pointer from one type of data structure to a pointer to a corresponding data structure, to sort a list of addresses, or to select the addresses of structures with certain properties.

MDB executes each dcmd in the pipeline in order from left to right. The left-most dcmd is executed using the current value of dot, or using the value specified by an explicit expression at the start of the command. When a | operator is encountered, MDB creates a pipe (a shared buffer) between the output of the dcmd to its left and the MDB parser, and an empty list of values.

As the dcmd executes, its standard output is placed in the pipe and then consumed and evaluated by the parser, as if MDB were reading this data from standard input. Each line must consist of an arithmetic expression terminated by a NEWLINE or semicolon (;). The value of the expression is appended to the list of values associated with the pipe. If a syntax error is detected, the pipeline is aborted.

When the dcmd to the left of a | operator completes, the list of values associated with the pipe is then used to invoke the dcmd to the right of the | operator. For each value in the list, dot is set to this value and the right-hand dcmd is executed. Only the rightmost dcmd in the pipeline has its output printed to standard output. If any dcmd in the pipeline produces output to standard error, these messages are printed directly to standard error and are not processed as part of the pipeline.

Formatting dcmds

The /, \setminus , ?, and = metacharacters are used to denote the special output formatting dcmds. Each of these dcmds accepts an argument list consisting of one or more format characters, repeat counts, or quoted strings. A format character is one of the ASCII characters shown in the table below.

Format characters are used to read and format data from the target. A repeat count is a positive integer preceding the format character that is always interpreted in base 10 (decimal). A repeat count can also be specified as an expression enclosed in square brackets preceded by a dollar sign (\$[]]). A string argument must be enclosed in double-quotes (" "). No blanks are necessary between format arguments.

The formatting dcmds are:

- / Display data from the target's virtual address space starting at the virtual address specified by dot.
- \ Display data from the target's physical address space starting at the physical address specified by dot.
- ? Display data from the target's primary object file starting at the object file location corresponding to the virtual address specified by dot.
- = Display the value of dot itself in each of the specified data formats. The = dcmd is therefore useful for converting between bases and performing arithmetic.

In addition to dot, MDB keeps track of another global value called the *increment*. The increment represents the distance between dot and the address following all the data read by the last formatting dcmd.

For example, if a formatting dcmd is executed with dot equal to address A, and displays a 4-byte integer, then after this dcmd completes, dot is still A, but the increment is set to 4. The + character, described in Arithmetic Expansion, would now evaluate to the value A + 4, and could be used to reset dot to the address of the next data object for a subsequent dcmd.

Most format characters increase the value of the increment by the number of bytes corresponding to the size of the data format, shown in the table. The table of format characters can be displayed from within MDB using the ::formats dcmd.

The format characters are:

- + Increment dot by the count (variable size)
- Decrement dot by the count (variable size)
- B Hexadecimal int (1 byte)
- C Character using C character notation (1 byte)
- D Decimal signed int (4 bytes)
- E Decimal unsigned long long (8 bytes)
- F Double (8 bytes)
- G Octal unsigned long long (8 bytes)
- H Swap bytes and shorts (4 bytes)
- I Address and disassembled instruction (variable size)
- J Hexadecimal long long (8 bytes)
- K Hexadecimal uintptr_t (4 or 8 bytes)
- N Newline
- O Octal unsigned int (4 bytes)

- P Symbol (4 or 8 bytes)
- Q Octal signed int (4 bytes)
- R Binary int (8 bytes)
- S String using C string notation (variable size)
- T Horizontal tab
- U Decimal unsigned int (4 bytes)
- V Decimal unsigned int (1 byte)
- W Default radix unsigned int (4 bytes)
- X Hexadecimal int (4 bytes)
- Y Decoded time32_t (4 bytes)
- Z Hexadecimal long long (8 bytes)
- ^ Decrement dot by increment * count (variable size)
- a Dot as symbol+offset
- b Octal unsigned int (1 byte)
- c Character (1 byte)
- d Decimal signed short (2 bytes)
- e Decimal signed long long (8 bytes)
- f Float (4 bytes)
- g Octal signed long long (8 bytes)
- h Swap bytes (2 bytes)
- i Disassembled instruction (variable size)
- n Newline
- o Octal unsigned short (2 bytes)
- p Symbol (4 or 8 bytes)
- q Octal signed short (2 bytes)
- r Whitespace
- s Raw string (variable size)
- t Horizontal tab
- u Decimal unsigned short (2 bytes)
- v Decimal signed int (1 byte)

- w Default radix unsigned short (2 bytes)
- x Hexadecimal short (2 bytes)
- y Decoded time64_t (8 bytes)

The /, \, and ? formatting dcmds can also be used to write to the target's virtual address space, physical address space, or object file by specifying one of the following modifiers as the first format character, and then specifying a list of words that are either immediate values or expressions enclosed in square brackets preceded by a dollar sign (\$[]).

The write modifiers are:

- v Write the lowest byte of the value of each expression to the target beginning at the location specified by dot
- w Write the lowest 2 bytes of the value of each expression to the target beginning at the location specified by dot
- W Write the lowest 4 bytes of the value of each expression to the target beginning at the location specified by dot
- Z Write the complete 8 bytes of the value of each expression to the target beginning at the location specified by dot

The /, \, and ? formatting dcmds can also be used to search for a particular integer value in the target's virtual address space, physical address space, and object file, respectively, by specifying one of the following modifiers as the first format character, then specifying a value and optional mask. The value and mask are each specified as either immediate values or expressions enclosed in square brackets preceded by a dollar sign.

If only a value is specified, MDB reads integers of the appropriate size and stops at the address containing the matching value. If a value V and mask M are specified, MDB reads integers of the appropriate size and stops at the address containing a value X where (X & M) = V. At the completion of the dcmd, dot is updated to the address containing the match. If no match is found, dot is left at the last address that was read.

The search modifiers are:

- 1 Search for the specified 2-byte value
- L Search for the specified 4-byte value
- M Search for the specified 8-byte value

For both user and kernel targets, an address space is typically composed of a set of discontiguous segments. It is not legal to read from an address that does not have a corresponding segment. If a search reaches a segment boundary without finding a match, it aborts when the read past the end of the segment boundary fails.

Chapter 4. Interaction

This chapter describes the MDB interactive command-line editing and history functions, the output pager, and debugger signal handling.

Command Re-entry

The text of the last HISTSIZE (default 128) commands entered from a terminal device are saved in memory. The in-line editing facility, described next, provides key mappings for searching and fetching elements from the history list.

In-line Editing

If standard input is a terminal device, MDB provides some simple emacs-style facilities for editing the command line. The search, previous, and next commands in edit mode provide access to the history list. Only strings, not patterns, are matched when searching. In the list below, the notation for control characters is caret (^) followed by a character shown in upper case. The notation for escape sequences is M- followed by a character. For example, M-f (pronounced *meta-eff*) is entered by depressing <ESC> followed by 'f', or by depressing Meta followed by 'f' on keyboards that support a Meta key. A command line is committed and executed using RETURN or NEWLINE. The edit commands are:

^F	Move cursor forward (right) one character.
M-f	Move cursor forward one word.
^B	Move cursor backward (left) one character.
M-b	Move cursor backward one word.
^A	Move cursor to start of line.
^E	Move cursor to end of line.
^D	Delete current character, if the current line is not empty. If the current line is empty, ^D denotes EOF and the debugger will exit.
M-^H	(Meta-backspace) Delete previous word.
^K	Delete from the cursor to the end of the line.
^ L	Reprint the current line.
^T	Transpose current character with next character.
^N	Fetch the next command from the history. Each time ^N is entered, the next command forward in time is retrieved.
^ P	Fetch the previous command from the history. Each time ^P is entered, the next command backward in time is retrieved.
^R[string]	Search backward in the history for a previous command line containing string. The string should be terminated by a RETURN or NEWLINE. If string is omitted, the previous

history element containing the most recent string is retrieved.

The editing mode also interprets the following user-defined sequences as editing commands. User defined sequences can be read or modified using the stty(1) command.

erase User defined erase character (usually ^H or ^?). Delete previous character.

intr User defined interrupt character (usually ^C). Abort the current command and print a new

prompt.

kill User defined kill character (usually ^U). Kill the entire current command line.

quit User defined quit character (usually ^\). Quit the debugger.

suspend User defined suspend character (usually ^Z). Suspend the debugger.

werase User defined word erase character (usually ^W). Erase the preceding word.

On keyboards that support an extended keypad with arrow keys, mdb will interpret these keystrokes as editing commands:

up-arrow Fetch the previous command from the history (same as ^P).

down-arrow Fetch the next command from the history (same as ^N).

left-arrow Move cursor backward one character (same as ^B).

right-arrow Move cursor forward one character (same as ^F)

Keyboard Shortcuts

MDB provides a set of keyboard shortcuts that bind individual keystrokes to common MDB commands when the keystroke listed in the table below is typed as the first character following the MDB prompt. The keyboard shortcuts are:

[Execute the command :: step over

] Execute the command ::step

Output Pager

mdb provides a built-in output pager. The output pager is enabled if the debugger's standard output is a terminal device. Each time a command is executed, mdb will pause after one screenful of output is produced and will display a pager prompt:

```
>> More [<space>, <cr>, q, n, c, a] ?
```

The following key sequences are recognized by the pager:

SPACE Display the next screenful of output.

a, A Abort the current top-level command and return to the prompt.

c, C Continue displaying output without pausing at each screenful until the

current top-level command is complete.

n, N, NEWLINE, RETURN Display the next line of output.

q, Q, ^C, ^\

Quit (abort) the current dcmd only.

Signal Handling

The debugger ignores the PIPE and QUIT signals. The INT signal aborts the command that is currently executing. The debugger intercepts and provides special handling for the ILL, TRAP, EMT, FPE, BUS, and SEGV signals. If any of these signals are generated asynchronously (that is, delivered from another process using kill(2)), mdb will restore the signal to its default disposition and dump core. However, if any of these signals are generated synchronously by the debugger process itself and a dcmd from an externally loaded dmod is currently executing, and standard input is a terminal, mdb will provide a menu of choices allowing the user to force a core dump, quit without producing a core dump, stop for attach by a debugger, or attempt to resume. The resume option will abort all active commands and unload the dmod whose dcmd was active at the time the fault occurred. It can then be subsequently re-loaded by the user. The resume option provides limited protection against buggy dcmds. Refer to Warnings, Use of the Error Recovery Mechanism, for information about the risks associated with the resume option.

Chapter 5. Built-in Commands

MDB provides a set of built-in dcmds that are always defined. Some of these dcmds are applicable only to certain targets: if a dcmd is not applicable to the current target, it fails and prints a message indicating "command is not supported by current target".

In many cases, MDB provides a mnemonic equivalent (::identifier) for the legacy adb(1) dcmd names. For example, ::quit is provided as the equivalent of \$q. Programmers who are experienced with adb(1) or who appreciate brevity or arcana might prefer the \$ or : forms of the built-ins. Programmers who are new to MDB might prefer the more verbose :: form. The built-ins are shown in alphabetical order. If a \$ or : form has a ::identifier equivalent, it is shown under the ::identifier form.

Built-in dcmds

> varAssignethe value of dot to the specified named variable. Some variables are read-only and cannot be name modified. If the > is followed by a modifier character surrounded by //, then the value is modified > /mods faire of the assignment. The modifier characters are:

```
variable-
namec Unsigned char quantity (1-byte)
```

- s Unsigned short quantity (2-byte)
- i Unsigned int quantity (4-byte)
- 1 Unsigned long quantity (4-byte in 32-bit, 8-byte in 64-bit)

Notice that these operators do not perform a cast; they instead fetch the specified number of low-order bytes (on little-endian architectures) or high-order bytes (big-endian architectures). These modifiers are provided for backward compatibility; the MDB */modifier/ and %/modifier/ syntax should be used instead.

\$< markened and execute commands from the specified macro file. The file name can be given as an absolute name or relative path. If the file name is a simple name (that is, if it does not contain a '/'), MDB searches for it in the macro file include path. If another macro file is currently being processed, this file is closed and replaced with the new file.

\$<< Read and execute commands from the specified macro file (as with \$<), but do not close the current *macropen* macro file.

name

]

\$? Print the process-ID and current signal of the target if it is a user process or core file, and then print the general register set of the representative thread.

[addPristsa C stack backtrace, including stack frame pointer information. If the dcmd is preceded by an explicit address, a backtrace beginning at this virtual memory address is displayed. Otherwise, the stack of the representative thread is displayed. If an optional count value is given as an argument, [cound more than count arguments are displayed for each stack frame in the output.

Note

The biased frame pointer value (that is, the virtual address minus 0x7ff) should be used as the address when requesting a stack trace.

[base et or set the default output radix. If the dcmd is preceded by an explicit expression, the default output radix is set to the given base; otherwise, the current radix is printed in base 10 (decimal).

- \$d The default radix is base 16 (hexadecimal).
- Print a list of all known external (global) symbols of type object or function, the value of the symbol, and the first 4 (32-bit **mdb**) or 8 (64-bit **mdb**) bytes stored at this location in the target's virtual address space. The **::nm** dcmd provides more flexible options for displaying symbol tables.

Pp Sentifice-prompt to the specified prompt-string. The default prompt is '>'. The prompt can also string set using ::set P or the P command-line option.

\$M In **kmdb** only, list the macro files that are cached by **kmdb** for use with the \$< dcmd.

distance for address-to-symbol-name conversions. The symbol matching distance modes are discussed along with the s command-line option in Appendix A, Options. The symbol matching distance can also be modified using the ::set s option. If no distance is specified, the current setting is displayed.

Print a list of the named variables that have non-zero values. The **::vars** dcmd provides other options for listing variables.

widt Set the output page width to the specified value. Typically, this command is not necessary, as MDB sw queries the terminal for its width and handles resize events.

Reopen the target for writing, as if MDB had been executed with the w option on the command line. Write mode can also be enabled with the ::set w option.

::arralyrint the address of each element of an array. The type of the array elements should be specified as typethe first argument, type, and the number of elements to be computed should be specified as the countecond argument, count. The output of ::array can be pipelined to the ::print dcmd to print the elements of an array data structure.

Note

This dcmd may only be used with objects that contain compressed symbolic debugging information designed for use with mdb. This information is currently only available for certain Solaris kernel modules. The SUNWzlib decompression software must be installed in order to process the symbolic debugging information.

[pidIf the user process target is active, attach to and debug the specified process-ID or core file. The]::attache file path name should be specified as a string argument. The process-ID can be specified as the [coretring argument, or as the value of the expression preceding the dcmd. Recall that the default base | pid is hexadecimal, so decimal PIDs obtained using pgrep(1) or ps(1) should be preceded with "0t"], [piwhen specified as expressions.

]
:A
[core
|pid
]

::brandheday the most recent branches taken by the current CPU. This dcmd is currently only available

[v when using kmdb on x86 systems where the appropriate processor-specific feature is enabled. The

number and type of branches that can be displayed is determined by the processor architecture. If
the v option is present, the instructions prior to each branch are displayed.

::cat Concatenate and display files. Each file name can be specified as a relative or absolute path name. file Thanfile contents will print to standard output, but will not pass through the output pager. This dcmd ... is intended to be used with the | operator; the programmer can initiate a pipeline using a list of addresses stored in an external file.

addresstext switch to the specified process. A context switch operation is valid only when using the ::contextnel target. The process context is specified using the address of its proc structure in the addressel's virtual address space. The special context address "0" is used to denote the context of the kernel itself. MDB can only perform a context switch when examining a crash dump if the dump contains the physical memory pages of the specified user process (as opposed to just kernel pages). The kernel crash dump facility can be configured to dump all pages or the pages of the current user process using dumpadm(1M). The ::status dcmd can be used to display the contents of the current crash dump.

When the user requests a context switch from the kernel target, MDB constructs a new target representing the specified user process. After the switch occurs, the new target interposes its dcmds at the global level: thus the / dcmd can now format and display data from the virtual address space of the user process, the ::mappings dcmd can display the mappings in the address space of the user process, and so on. The kernel target can be restored by executing 0::context.

::cpuidigsplay the current general-purpose register set for the current CPU or the specified cpuid. This [c cpaonamand is only available when using kmdb.

::cpustbacklay a C stack backtrace for the thread executing on the current CPU or the specified cpuid. [c cpthistcommand is only available when using kmdb.]

::dcmlisst the available dcmds and print a brief description for each one.

[addDissassemble starting at or around the address specified by the final argument, or the current value] ::disof dot. If the address matches the start of a known function, the entire function is disassembled; afwotherwise, a "window" of instructions before and after the specified address is printed in order to provide context. By default, instructions are read from the target's virtual address space; if the f [n contion is present, instructions are read from the target's object file instead. The f option is enabled] [addit default if the debugger is not currently attached to a live process, core file, or crash dump. The w option can be used to force window-mode, even if the address is the start of a known function. The size of the window defaults to ten instructions; the number of instructions can be specified explicitly using the n option. If the a option is present, addresses are printed as numeric values rather than symbolically.

::disakinks the available disassembler modes. When a target is initialized, MDB attempts to select the appropriate disassembler mode. The user can change the mode to any of the modes listed using the ::dismode dcmd.

::disn@decor set the disassembler mode. If no argument is specified, print the current disassembler mode. [mode argument is specified, switch the disassembler to the specified mode. The list of available], disassemblers can be displayed using the ::disasms dcmd.

\$V [mode]

::dmoldist the loaded debugger modules. If the 1 option is specified, the list of the dcmds and walkers
 [1] associated with each dmod is printed below its name. The output can be restricted to a particular dmod by specifying its name as an additional argument.

```
][module-
name
]
```

[address] hexadecimal and ASCII memory dump of the 16-byte aligned region of virtual memory ::dumpontaining the address specified by dot. If a repeat count is specified for ::dump, this is interpreted [eqrstu]number of bytes to dump rather than a number of iterations. The ::dump dcmd also recognizes [f|p] the following options:

[g bytes]

[w pa£agraphs]	Adjust for endian-ness. The e option assumes 4-byte words; the g option can be used to change the default word size.	
f	Read data from the object file location corresponding to the given virtual address instead of from the target's virtual address space. The f option is enabled by default if the debugger is not currently attached to a live process, core file, or crash dump.	
g group	Display bytes in groups of bytes. The default <i>group</i> size is 4 bytes. The <i>group</i> size must be a power of two that divides the line width	
р	Interpret address as a physical address location in the target's address space instead of a virtual address.	
đ	Do not print an ASCII decoding of the data.	
r	Number lines relative to the start address instead of with the explicit address of each line. This option implies the u option.	
s	Elide repeated lines.	
t	Only read from and display the contents of the specified addresses, instead of reading and printing entire lines.	
u	Unalign output instead of aligning the output at a paragraph boundary.	
w paragraphs	Display paragraphs 16-byte paragraphs per line. The default number of paragraphs is one. The maximum value accepted for w is 16.	

::echdPrint the arguments separated by blanks and terminated by a NEWLINE to standard output. [strExpressions enclosed in \$[] will be evaluated to a value and printed in the default base. |value

...]

::evalEvaluate and execute the specified string as a command. If the command contains metacharacters command it space, it should be enclosed in double or single quotes.

::filesPrint a list of the known source files (symbols of type STT_FILE present in the various target [objesymbol tables). If an object name is specified, the output is restricted to file symbols present in the corresponding object file.

[add Sesss] instruction text for instructions that refer to the specified symbols or addresses. The search ::findsymshould consist of one or more addresses or symbol names specified as an address preceding the [g] dcmd or one or more symbol names or expressions following the dcmd. If the g option is specified, [additional additional additi

...]

Note

SPARC only. The ::findsym dcmd is only available when debugging a target that uses the SPARC instruction set architecture.

::formlats the available output format characters for use with the /, \, ?, and = formatting dcmds. The formats and their use is described in Formatting dcmds.

[thread in the floating-point register set of the representative thread. If a thread is specified, the floating]::fpressnt registers of that thread are displayed. The thread expression should be one of the thread [dqs]identifiers described under Thread Support.

```
, [ thread ]

Note
$x, $X,
$y,
SPARC only. The d, q, and s options can be used to display the floating point registers as a collection of double-precision (d), quad-precision (q), or single-precision (s) floating point values.
```

::grepEvaluate the specified command string, then print the old value of dot if the new value of dot is commandzero. If the command contains white space or metacharacters, it must be quoted. The **::grep** dcmd can be used in pipelines to filter a list of addresses.

::helpWith no arguments, the **::help** dcmd prints a brief overview of the help facilities available in MDB. [dcmld-a dcmd-name is specified, MDB prints a usage summary for that dcmd. name

[addReads and display len bytes from the I/O port specified by address. The value of the L option, [, leif present, takes precedence over the repeat count specified on the left-hand side. The len must be]]::in1, 2, or 4 bytes and the port address must be aligned according to the length. This command is only [L leavailable when using kmdb on x86 systems.]

[add Wests through the elements of a linked list data structure and print the address of each element in] ::list the list. The address of the first element in the list can be specified using an optional address; typeotherwise the list is assumed to start at the current value of dot. The type parameter must name member struct or union type and is used to describe the type of the list elements so that MDB can read [varinable ests of the appropriate size. The member parameter is used to name the member of type namethat contains a pointer to the next list element. The ::list dcmd will continue iterating until a NULL pointer is encountered, the first element is reached again (a circular list), or an error occurs while reading an element. If the optional variable-name is specified, the specified variable will be assigned the value returned at each step of the walk when MDB invokes the next stage of a pipeline.

Note

1

This dcmd may only be used with objects that contain compressed symbolic debugging information designed for use with mdb. This information is currently only available for certain Solaris kernel modules. The SUNWzlib decompression software must be installed in order to process the symbolic debugging information.

::loadLoad the specified dmod. The module name can be given as an absolute or relative path. If modulename is a simple name (that is, does not contain a '/'), MDB searches for it in the module library
s] modules with conflicting names cannot be loaded; the existing module must be unloaded first.
nameIf the -s option is present, MDB will remain silent and not issue any error messages if the module is not found or could not be loaded.

::log Enable or disable the output log. MDB provides an interactive logging facility where both the input [d commands and standard output can be logged to a file while still interacting with the user. The e option enables logging to the specified file, or re-enables logging to the previous log file if no file e] f inhance is given. The d option disables logging. If the \$> dcmd is used, logging is enabled if a file], name argument is specified; otherwise, logging is disabled. If the specified log file already exists, \$> MDB appends any new log output to the file.

[filename]

::mapMap the value of dot to a corresponding value using the command specified as a string argument, commthed print the new value of dot. If the command contains white space or metacharacters, it must be quoted. The ::map dcmd can be used in pipelines to transform the list of addresses into a new list of addresses.

[addPristsa list of each mapping in the target's virtual address space, including the address, size, and]::mappingstion of each mapping. If the dcmd is preceded by an address, MDB shows only the mapping [namethat contains the given address. If a string name argument is given, MDB shows only the mapping], that matched the description.

```
[ address
]
$m
[name
```

[address] he symbol tables associated with the current target. If an optional address preceding the ::nm dcmd is specified, only the symbol table entry for the symbol corresponding to address is displayed. [DPdfframophy.cct name is specified, only the symbol table for this load object is displayed. The ::nm] [dcmd also recognizes the following options:

```
t types
][ D
f format
][object]
d
```

Prints . dynsym (dynamic symbol table) instead of . symtab.

Prints the private symbol table instead of .symtab.

Prints value and size fields in decimal.

f format [,format...]

Print only the specified symbol information. The valid format argument strings are:

```
ndx
         symbol table index
val
         symbol table
size
         size in bytes
         symbol type
type
bind
         binding
oth
         other
         section index
shndx
         symbol name
name
         C type for symbol (if known)
ctype
         object which defines symbol
obj
```

g	Prints o	only global symbols.
h	Suppre	sses the header line.
n	Sorts sy	mbols by name.
0	Prints v	ralue and size fields in octal.
р	be used	ymbols as a series of ::nmadd commands. This option can with P to produce a macro file that can be subsequently o the debugger with \$<.
t type [,type]		only symbols of the specified type(s). The valid type nt strings are:
	noty	STT_NOTYPE
	objt	STT_OBJECT
	func	STT_FUNC
	sect	STT_SECTION
	file	STT_FILE
	comm	STT_COMMON
	tls	STT_TLS
	regi	STT_SPARC_REGISTER
u	Prints o	only undefined symbols.
v	Sorts sy	mbols by value.
x	Prints v	alue and size fields in hexadecimal.
	nterpose	ivate symbol table. MDB provides a private, configurable on the target's symbol table, as described in Symbol Name gnizes the following options:

Resolution. The **::nmadd** dcmd also recognizes the following options:

fo

Set the size of the symbol to end - value. е

[eend

f Set the type of the symbol to STT_FUNC.

[ssize

] name Set the type of the symbol to STT_OBJECT.

Set the size of the symbol to size.

::nmdeelete the specified symbol name from the private symbol table. name

::objettsnt a map of the target's virtual address space, showing only those mappings that correspond to the primary mapping (usually the text section) of each of the known load objects. If the v option is present, the command displays the version of each object if version information is know. If no version information is known, a version of Unknown will be displayed in the output.

::offsetnit the offset of the specified member of the specified type. The type should be the name of a C typestructure. The offset is printed in bytes, unless the member is a bit-field in which case the offset may member printed in bits. The output is always suffixed with the appropriate units for clarity. The type name may use the backquote (`) scoping operator described in Symbol Name Resolution.

Note

This dcmd may only be used with objects that contain compressed symbolic debugging information designed for use with mdb. This information is currently only available for certain Solaris kernel modules. The SUNWzlib decompression software must be installed in order to process the symbolic debugging information.

[add Write the specified value to the I/O port specified by address. The value of the L option, if [, lepresent, takes precedence over the repeat count specified on the left-hand side. The len must be 1,]] ::out, or 4 bytes and the port address must be aligned according to the length. This command is only [L lenvailable when using kmdb on x86 systems.]

[address] the data structure at the specified virtual address using the given type information. The ::printype parameter may name a C struct, union, enum, fundamental integer type, or a pointer to any [acdiffical] types. If the type name contains whitespace (for example, "struct foo") it must be enclosed [c lim] single or double quotes. The type name may use the backquote (`) scoping operator described [lim] der Symbol Name Resolution. If the type is a structured type, the ::print dcmd will recursively [typeprint each member of the struct or union. If the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type argument is not present and a static or global [member of the struct or union the type

The type argument can be followed by an optional list of member or offset expressions, in which case only those members and submembers of the specified type are displayed. Members can be specified using C syntax that includes the array index operator ([]), the structure member operator (->), and the structure pointer operator (.). Offsets can be specified using the MDB arithmetic expansion syntax (\$[]). After displaying the data structure, ::print increments dot by the size of type in bytes.

Note

The ::print dcmd may only be used with objects that contain compressed symbolic debugging information designed for use with MDB. This information is only available at present in certain Solaris kernel modules and user libraries. The SUNWzlib decompression software must be installed in order to process the symbolic debugging information.

If the a option is present, the address of each member is displayed. If the i option is present, the expression on the left-hand side is interpreted as an immediate value to be displayed using the specified type. If the p option is present, ::print interprets address as a physical memory address instead of a virtual memory address. If the t option is present, the type of each member is displayed. If the d or x options are present, all integers are displayed in decimal (d) or hexadecimal (x); by default a heuristic is used to determine if the value should be displayed in decimal or hexadecimal. The number of characters in a character array that will be read and displayed as a string can be limited with the c option. If the C option is present, no limit is enforced. The number of elements in a standard array that will be read and displayed can be limited with the 1 option. If the L option is present, no limit is enforced and all array elements are shown. The default values for c and 1 can be modified using : set or the o command-line option as described in Appendix A, Options.

::quitQuit the debugger. When using **kmdb** only, the u option causes the debugger to resume execution of [u the operating system and unload the debugger. The u option cannot be used if **kmdb** was loaded at],

```
boot. If the u option is not present, ::quit causes kmdb to exit to the firmware (on SPARC systems)
[ u or causes the system to reboot (on x86 systems).
1
[ thr Prior the general-purpose register set of the representative thread. If a thread is specified, the general
::regsurpose register set of that thread is displayed. The thread expression should be one of the thread
[ thricked tifiers described under Thread Support.
1
$r
::released ease the previously attached process or core file. If the -a option is present, the process is released
     and left stopped and abandoned. It can subsequently be continued by prun (1) or it can be resumed
     by applying MDB or another debugger. By default, a released process is forcibly terminated if it was
1,
     created by MDB using ::run, or it is released and set running if it was attached to by MDB using
     the -p option or using the ::attach or :A dcmds.
1
::set Get or set miscellaneous debugger properties. If no options are specified, the current set of debugger
      properties is displayed. The ::set dcmd recognizes the following options:
ſ
wF
         Forcibly take over the next user process that ::attach is applied to, as if mdb had been executed
11
          with the F option on the command line.
+/
-o option
      I Set the default path for locating macro files. The path argument can contain any of the special
[s distanteens described for the I command-line option in Appendix A, Options.
I path Set the default path for locating debugger modules. The path argument can contain any of the
          special tokens described for the I command-line option in Appendix A, Options.
1
```

- P Set the command prompt to the specified prompt string.
- s Set the symbol matching distance to the specified distance. Refer to the description of the s command-line option in Appendix A, Options for more information.

o Enable the specified debugger option. If the +o form is used, the option is disabled. The option

[P prompttings are described along with the o command-line option in Appendix A, Options.

w Re-open the target for writing, as if mdb had been executed with the w option on the command line.

::showDisplay revision information for the hardware and software corresponding the current target. If no options are specified, general system information is displayed. If the p option is present, information pv for each load object that is part of a patch is displayed. If the v option is present, information for each load object is displayed. Load objects without version information will be omitted from the output for the p option. Load objects without version information will report Unknown in the output of the v option.

::sizeoffrint the size of the specified type in bytes. The *type* parameter may name a C struct, union, enum, *type*fundamental integer type, or a pointer to any of these types. The type name may use the backquote (`) scoping operator described in Symbol Name Resolution.

Note

[Lpath]

This dcmd may only be used with objects that contain compressed symbolic debugging information designed for use with mdb. This information is currently only available for

certain Solaris kernel modules. The SUNWzlib decompression software must be installed in order to process the symbolic debugging information.

[addPristsa C stack back trace. If the dcmd is preceded by an explicit address, a back trace beginning] ::stackthis virtual memory address is displayed. Otherwise, the stack of the representative thread is [courtisplayed. If an optional count value is given as an argument, no more than count arguments are], displayed for each stack frame in the output.

[address | Note

\$c

[count The biased frame pointer value (that is, the virtual address minus 0x7ff) should be used as the address when requesting a stack trace.

::staturint a summary of information related to the current target.

cpuil When using **kmdb** only, switch to the CPU indicated by the specified cpuid and use this CPU's **::switch**; rent register state as the representative for debugging.

:X

::ternPrint the name of the terminal type that MDB is using to perform any terminal-dependent input and output operations, such as command-line editing.

thre Reint the address of the storage for the specified thread-local storage (TLS) symbol in the context of ::tls the specified thread. The thread expression should be one of the thread identifiers described under symboliread Support. The symbol name may use any of the scoping operators described under Symbol Name Resolution.

::typeSet attributes for named variables. If one or more variable names are specified, they are defined and [+/ set to the value of dot. If the t option is present, the user-defined tag associated with each variable -t] variable to option is present, the tag is cleared. If no variable names are specified, the list of name variables and their values is printed.

...

::unloadload the specified dmod. The list of active dmods can be printed using the ::dmods dcmd. Built-moduihemodules cannot be unloaded. Modules that are busy (that is, provide dcmds that are currently name executing) cannot be unloaded.

::unsetInset (remove) the specified variable(s) from the list of defined variables. Some variables are varied by MDB are marked as persistent, and cannot be unset by the user.

name

...

::varsPrint a listing of named variables. If the n option is present, the output is restricted to variables that [npt]currently have non-zero values. If the p option is present, the variables are printed in a form suitable for re-processing by the debugger using the \$< dcmd. This option can be used to record the variables to a macro file, then restore these values later. If the t option is present, only the tagged variables are printed. Variables can be tagged using the t option of the ::typeset dcmd.

::version the debugger version number.

address the physical address mapping for the specified virtual address, if possible. The ::vtop dcmd ::vtop only available when examining a kernel target, or when examining a user process inside a kernel [a astrash dump (after a ::context dcmd has been issued).

When examining a kernel target from the kernel context, the a option can be used to specify the address (as) of an alternate address space structure that should be used for the virtual to physical

translation. By default, the kernel's address space is used for translation. This option is available for active address spaces even when the dump content only contains kernel pages.

[add Washsthrough the elements of a data structure using the specified walker. The available walkers can] ::walke listed using the ::walkers dcmd. Some walkers operate on a global data structure and do not walkerquire a starting address. For example, walk the list of proc structures in the kernel. Other walkers nameoperate on a specific data structure whose address must be specified explicitly. For example, given [vari pointer to an address space, walk the list of segments.

name

When used interactively, the ::walk dcmd will print the address of each element of the data structure in the default base. The dcmd can also be used to provide a list of addresses for a pipeline. The walker name can use the backquote " " " scoping operator described in dcmd and Walker Name Resolution. If the optional variable-name is specified, the specified variable will be assigned the value returned at each step of the walk when MDB invokes the next stage of the pipeline.

:: walkers the available walkers and print a brief description for each one.

::whereint the dmod that exports the specified dcmds and walkers. These dcmds can be used to determine [v] nawheich dmod is currently providing the global definition of the given dcmd or walker. Refer to dcmd ..., ::white Walker Name Resolution for more information on global name resolution. The v option causes [v] nathedcmd to print the alternate definitions of each dcmd and walker in order of precedence.

::xdataist the external data buffers exported by the current target. External data buffers represent information associated with the target that cannot be accessed through standard target facilities (that is, an address space, symbol table, or register set). These buffers can be consumed by dcmds; for more information, refer to mdb_get_xdata().

Chapter 6. Execution Control

MDB provides facilities for controlling and tracing the execution of live running programs, including both user applications and the live operating system kernel and device drivers. You can use the **mdb** command to control user processes that are already running, or create new processes under the control of the debugger. You can boot or load **kmdb** to control the execution of the operating system kernel itself, or debug a device driver. This chapter describes the built-in dcmds that can be used to control target execution. These commands can be used in either **mdb** or **kmdb**, except as noted in the descriptions. Additional topics relating only to execution control in **kmdb** are discussed in Chapter 7, Kernel Execution Control.

Execution Control

MDB provides a simple model of execution control: a target process can be started from within the debugger using ::run, or MDB can attach to an existing process using :A, ::attach, or the -p command-line option (see Chapter 5, Built-in Commands). Alternately, the kernel can be booted using kmdb or kmdb can be loaded afterward. In either case, a list of traced software events can be specified by the user. Each time a traced event occurs in the target program, all threads in the target stop, the thread that triggered the event is chosen as the representative thread, and control returns to the debugger. Once the target program is set running, control can be asynchronously returned to the debugger by typing the user-defined interrupt character (typically Control-C).

A *software event* is a state transition in the target program that is observed by the debugger. For example, the debugger may observe the transition of a program counter register to a value of interest (a breakpoint) or the delivery of a particular signal.

A software event specifier is a description of a class of software events that is used by the debugger to instrument the target program in order to observe these events. The ::events dcmd is used to list the software event specifiers. A set of standard properties is associated with each event specifier, as described under ::events in Built-in dcmds.

The debugger can observe a variety of different software events, including breakpoints, watchpoints, signals, machine faults, and system calls. New specifiers can be created using ::bp, ::fltbp, :: sigbp, ::sysbp, or ::wp. Each specifier has an associated callback (an MDB command string to execute as if it had been typed at the command prompt) and a set of properties, as described under ::events in Built-in dcmds. Any number of specifiers for the same event may be created, each with different callbacks and properties. The current list of traced events and the properties of the corresponding event specifiers can be displayed using the ::events dcmd. The event specifier properties are defined as part of the description of the ::events and ::evset dcmds, in Built-in dcmds.

The execution control built-in dcmds, described in Built-in dcmds, are always available, but will issue an error message indicating they are not supported if applied to a target that does not support execution control.

Event Callbacks

The **::evset** dcmd and event tracing dcmds allow you to associate an event callback (using the -c option) with each event specifier. The event callbacks are strings that represent MDB commands to execute when the corresponding event occurs in the target. These commands are executed as if they had been typed at the command prompt. Prior to executing each callback, the *dot* variable is set to the value of the representative thread's program counter and the *hits* variable is set to the number of times this specifier has been matched, including the current match.

If the event callbacks themselves contain one or more commands to continue the target (for example, ::cont or ::step), these commands do *not* immediately continue the target and wait for it to stop again. Instead, inside of an event callback, the continue dcmds note that a continue operation is now pending, and then return immediately. Therefore, if multiple dcmds are included in an event callback, the step or continue dcmd should be the last command specified. Following the execution of *all* event callbacks, the target will immediately resume execution if *all* matching event callbacks requested a continue. If conflicting continue operations are requested, the operation with the highest precedence determines what type of continue will occur. The order of precedence from highest to lowest is: step, step-over (next), step-out, continue.

Thread Support

MDB provides facilities to examine the stacks and registers of each thread associated with the target. The persistent "thread" variable contains the current representative thread identifier. The format of the thread identifier depends on the target. The ::regs and ::fpregs dcmds can be used to examine the register set of the representative thread, or of another thread if its register set is currently available. In addition, the register set of the representative thread is exported as a set of named variables. The user can modify the value of one or more registers by applying the > dcmd to the corresponding named variable.

The MDB kernel target exports the virtual address of the corresponding internal thread structure as the identifier for a given thread. This address corresponds to the kthread_t data structure in the operating system source code. When using **kmdb**, the CPU identifier for the CPU running **kmdb** is stored in the cpuid variable.

The MDB process target provides proper support for examination of multi-threaded user processes that use the native lwp_* interfaces, /usr/lib/libthread.so, or /usr/lib/libpthread.so. When debugging a live user process, MDB will detect if a single threaded process dlopens or closes libthread and will automatically adjust its view of the threading model on-the-fly. The process target thread identifiers will correspond to either the lwpid_t, thread_t, or pthread_t of the representative, depending on the threading model used by the application.

If MDB is debugging a user process target and the target makes use of compiler-supported thread-local storage, MDB will automatically evaluate symbol names referring to thread-local storage to the address of the storage corresponding to the current representative thread. The ::tls built-in dcmd can be used to display the value of the symbol for threads other than the representative thread.

Built-in dcmds

[addSet a breakpoint at the specified locations. The ::bp dcmd sets a breakpoint at each address or symbol ::bpspecified, including an optional address specified by an explicit expression preceding the dcmd, and each string or immediate value following the dcmd. The arguments may either be symbol names or -dDestimediate values denoting a particular virtual address of interest. If a symbol name is specified, it [c cmonlay refer to a symbol that cannot yet be evaluated in the target process: that is, it may consist of In coamobject name and function name in a load object that has not yet been opened. In this case, the sym breakpoint is deferred and it will not be active in the target until an object matching the given name ..., adistroaded. The breakpoint will be automatically enabled when the load object is opened. Breakpoints :b [crown symbols defined in a shared library should always be set using a symbol name and not using an address expression, as the address may refer to the corresponding Procedure Linkage Table (PLT) entry instead of the actual symbol definition. Breakpoints set on PLT entries may be overwritten by the run-time link-editor when the PLT entry is subsequently resolved to the actual symbol definition. The d, D, e, s, t, T, c, and n options have the same meaning as they do for the ::evset dcmd, as described later in this section. If the :b form of the dcmd is used, a breakpoint is only set at the virtual address specified by the expression preceding the dcmd. The arguments following the :b demd are concatenated together to form the callback string. If this string contains meta-characters, it must be quoted.

function using **kmdb** only, call the specified *function* defined in the operating system kernel. The **::call** *function* expression must match the address of a defined function in a symbol table of one of the [**arg**known kernel modules. If expression arguments are specified, these arguments as passed by value. ...] If string arguments are specified, these arguments are passed by reference.

Note

The ::call command should be used only with extreme caution and should never be applied to a production system. The operating system kernel will not resume execution in order to execute the specified function. Therefore, the function being called must not utilize arbitrary kernel services and must not block for any reason. You must be fully aware of the side-effects of any function you call using this command.

::contSuspend the debugger, continue the target program, and wait for it to terminate or stop following [SIG]a software event of interest. If the target is already running because the debugger was attached to :c [SIG]unning program with the -o nostop option enabled, this dcmd simply waits for the target to terminate or stop after an event of interest. If an optional signal name or number (see the signal (3HEAD) man page) is specified as an argument, the signal is immediately delivered to the target as part of resuming its execution. If the SIGINT signal is traced, control may be asynchronously returned to the debugger by typing the user-defined interrupt character (usually ^C). This SIGINT signal will be automatically cleared and will not be observed by the target the next time it is continued. If no target program is currently running, ::cont will start a new program running as if by ::run.

addrDelete the event specifiers with the given id number. The id number argument is interpreted in ::deletecimal by default. If an optional address is specified preceding the dcmd, all event specifiers that [id are associated with the given virtual address are deleted (e.g. all breakpoints or watchpoints affecting | all]that address). If the special argument "all" is given, all event specifiers are deleted, except those addrthat are marked sticky (T flag). The ::events dcmd displays the current list of event specifiers.
id [id

|all

::evDisptay the list of software event specifiers. Each event specifier is assigned a unique ID number [av], that can be used to delete or modify it at a later time. The debugger may also have its own internal **\$b** [av] ents enabled for tracing; these will only be displayed if the a option is present. If the v option is present, a more verbose display including the reason for any specifier inactivity will be shown. The following ::events dcmd shows example output:

```
> ::events
  ID S TA HT LM Description
                                                  Action
  ___ _ __ __ __
[ 1 ] - T
          1
            0 stop on SIGINT
[ 2 ] - T
          0
            0 stop on SIGQUIT
[ 3 ] - T
            0 stop on SIGILL
            0 stop on SIGXCPU
[ 11] - T
          0
[12] - T
            0 stop on SIGXFSZ
          2
            0 stop at libc`printf
[ 13] -
                                                  ::echo printf
```

The following discussion explains the meaning of each column. A summary of this information is available using ::help events.

ID The event specifier identifier. The identifier will be shown in square brackets [] if the specifier is enabled, in parentheses () if the specifier is disabled, or in angle

brackets <> if the target program is currently stopped on an event that matches the given specifier.

S	The event specifier state.	The state will be	one of the following syr	nbols:

-	The event specifier is idle. When no target program is running, all specifiers
	are idle. When the target program is running, a specifier may be idle if it
	cannot be evaluated (such as a deferred breakpoint in a shared object that
	is not yet loaded).

- + The event specifier is active. When the target is continued, events of this type will be detected by the debugger.
- * The event specifier is armed. This state means that the target is currently running with instrumentation for this type of event. This state is only visible if the debugger is attached to a running program with theo nostop option.
- ! The event specifier was not armed due to an operating system error. The ::events -v option can be used to display more information about the reason the instrumentation failed.

TA The Temporary, Sticky, and Automatic event specifier properties. One or more of the following symbols may be shown:

t	The event specifier is temporary, and will be deleted the next time the	
	stops, regardless of whether it is matched.	

- The event specifier is sticky, and will be not be deleted by ::delete all or :z. The specifier can be deleted by explicitly specifying its id number to::delete.
- d The event specifier will be automatically disabled when the hit count is equal to the hit limit.
- D The event specifier will be automatically deleted when the hit count is equal to the hit limit.
- s The target will automatically stop when the hit count is equal to the hit limit.

HT The current hit count. This column displays the number of times the corresponding software event has occurred in the target since the creation of this event specifier.

The current hit limit. This column displays the limit on the hit count at which the auto-disable, auto-delete, or auto-stop behavior will take effect. These behaviors can be configured using the ::evset dcmd.

Description A description of the type of software event that is matched by the given specifier.

The callback string to execute when the corresponding software event occurs. This callback is executed as if it had been typed at the command prompt.

id::eVsedify the properties of one or more software event specifiers. The properties are set for each [+/ specifier identified by the optional expression preceding the dcmd and an optional list of arguments -dDesfollowing the dcmd. The argument list is interpreted as a list of decimal integers, unless an explicit [c cmcdkix is specified. The ::evset dcmd recognizes the following options: [n count]

Disable the event specifier when the hit count reaches the hit limit. If the +d form of the option is given, this behavior is disabled. Once an event specifier is disabled, the debugger will remove

LM

Action

- any corresponding instrumentation and will ignore the corresponding software events until the specifier is subsequently re-enabled. If the n option is not present, the specifier is disabled immediately.
- Delete the event specifier when the hit count reaches the hit limit. If the +D form of the option is given, this behavior is disabled. The D option takes precedence over the d option. The hit limit can be configured using the n option.
- e Enable the event specifier. If the +e form of the option is given, the specifier is disabled.
- s Stop the target program when the hit count reaches the hit limit. If the +s form of the option is given, this behavior is disabled. The s behavior tells the debugger to act as if ::cont were issued following each execution of the specifier's callback, except for the Nth execution, where N is the current value of the specifier's hit limit. The s option takes precedence over both the D option and the d option.
- -t Mark the event specifier as temporary. Temporary specifiers are automatically deleted the next time the target stops, regardless of whether it stopped as the result of a software event corresponding to the given specifier. If the +t form of the option is given, the temporary marker is removed. The -t option takes precedence over the -T option.
- Mark the event specifier as sticky. Sticky specifiers will not be deleted by ::delete all or :z. They can be deleted by specifying the corresponding specifier ID as an explicit argument to ::delete. If the +T form of the option is given, the sticky property is removed. The default set of event specifiers are all initially marked sticky.
- Execute the specified *cmd* string each time the corresponding software event occurs in the target program. The current callback string can be displayed using ::events.
- n Set the current value of the hit limit to *count*. If no hit limit is currently set and the n option does not accompany s or D, the hit limit will be set to one.

A summary of this information is available using ::help evset.

*:fltbphe dcmd, or a list of fault names or numbers (see <sys/fault.h>) following the dcmd. The d, [+/ D, e, s, t, T, c, and n options have the same meaning as they do for the ::evset dcmd. The ::fltbp-dDestbhmand applies to user process debugging only.

```
[c cmd]
[n count]
flt
```

sign#the target is a live user process, ignore the specified signal and allow it to be delivered transparently

- to the target. All event specifiers that are tracing delivery of the specified signal will be deleted from the list of traced events. By default, the set of ignored signals is initialized to the complement of the set of signals that cause a process to dump core by default (see the signal(3HEAD) man page), except for SIGINT, which is traced by default. The :i command applies to user process debugging only.
- \$i Display the list of signals that are ignored by the debugger and will be handled directly by the target. More information on traced signals can be obtained using the ::events dcmd. The \$i command applies to user process debugging only.

::kill, Forcibly terminate the target if it is a live user process. The target will also be forcibly terminated
:k when the debugger exits if it was created by the debugger using ::run. The ::kill command applies to user process debugging only.

- \$1 Print the LWPID of the representative thread, if the target is a user process.
- **\$L** Print the LWPIDs of each LWP in the target, if the target is a user process.

::nextStep the target program one instruction, but step over subroutine calls. If an optional signal name [SIG]or number (see signal(3HEAD) man page) is specified as an argument, the signal is immediately :e [SIG]ivered to the target as part of resuming its execution. If no target program is currently running, ::next will start a new program running as if by ::run and stop at the first instruction.

::run Start a new target program running with the specified arguments and attach to it. The arguments are [argsnot interpreted by the shell. If the debugger is already examining a live running program, it will first ... detach from this program as if by ::release.

```
],
:r [args
...]
```

[sigfladde delivery of the specified signals. The signals are identified using an optional signal number ::sigbpreceding the dcmd, or a list of signal names or numbers (see signal(3HEAD)) following the dcmd. [+/ The d, D, e, s, t, T, c, and n options have the same meaning as they do for the ::evset dcmd. -dDesflitially, the set of signals that cause the process to dump core by default (see signal(3HEAD)) and [c cms]GINT are traced. The ::sigbp command applies to user process debugging only.

```
[n count]
SIG
..., [signal]
:t
[+/dDestT]
[c cmd]
[n count]
SIG
...
```

::stepStep the target program one instruction. If an optional signal name or number (see the [brarsignal(3HEAD) man page) is specified as an argument and the target is a user process, the signal ove is immediately delivered to the target as part of resuming its execution. If the optional branch out argument is specified, the target program will continue until the next instruction that branches the [SIG]control flow of the processor. The ::step branch feature is only available when using kmdb on :s SIG.86 systems with appropriate processor-specific features enabled. If the optional over argument :u SIG specified, ::step will step over subroutine calls. The ::step over argument is the same as the ::next dcmd. If the optional out argument is specified, the target program will continue until the representative thread returns from the current function. If no target program is currently running, ::step over will start a new program running as if by ::run and stop at the first instruction. The :s dcmd is the same as ::step. The :u dcmd is the same as ::step out.

[sys&ade]entry to or exit from the specified system calls. The system calls are identified using an optional ::sysbpystem call number preceding the dcmd, or a list of system call names or numbers (see <sys/[+/dDeşbade]1.h>) following the dcmd. If the i option is specified (the default), the event specifiers [io] trigger on entry into the kernel for each system call. If the o option is specified, the event specifiers [c cmt]gger on exit out from the kernel. The d, D, e, s, t, T, c, and n options have the same meaning as [n cothet] do for the ::evset dcmd. The ::sysbp command applies to user process debugging only. syscall

...

addrSet a watchpoint at the specified address. The length in bytes of the watched region may be set by [,len]specifying an optional repeat count preceding the dcmd. If no length is explicitly set, the default is

[+/dDesstb]te. The ::wp dcmd allows the watchpoint to be configured to trigger on any combination of [rwx]read (r option), write (w option), or execute (x option) access. The d, D, e, s, t, T, c, and n options [ip] have the same meaning as they do for the ::evset dcmd. When using kmdb on x86 systems only, [c cmt]e i option can be used to indicate that a watchpoint should be set on the address of an I/O port. In colling kmdb only, the p option can be used to indicate that the specified address should be addrinterpreted as a physical address. The :a dcmd sets a read access watchpoint at the specified address. [,len]the :p dcmd sets an execute access watchpoint at the specified address. The :w dcmd sets a write [cmd.access watchpoint at the specified address. The arguments following the :a. :p, and :w dcmds are], addxoncatenated together to form the callback string. If this string contains meta-characters, it must be [,len]the the callback string is the configuration of the callback string contains meta-characters, it must be [,len]the the callback string is the configuration of the callback string contains meta-characters, it must be [,len]the the callback string is the configuration of the callback string contains meta-characters, it must be [,len]the callback string is the configuration of the callback string is the cal

```
[cmd...], addr
[,len]:w
[cmd...]
```

Delete all event specifiers from the list of traced software events. Event specifiers can also be deleted using ::delete.

Interaction With exec

When a controlled user process performs a successful **exec**(2), the behavior of the debugger is controlled by the ::set -o follow_exec_mode option, as described in Summary of Command-line Options. If the debugger and victim process have the same data model, then the "stop" and "follow" modes determine whether MDB automatically continues the target or returns to the debugger prompt following the **exec**. If the debugger and victim process have a different data model, then the "follow" behavior causes MDB to automatically re-exec the MDB binary with the appropriate data model and re-attach to the process, still stopped on return from the **exec**. Not all debugger state is preserved across this re-exec.

If a 32-bit victim process **execs** a 64-bit program, then "stop" will return to the command prompt, but the debugger will no longer be able to examine the process because it is now using the 64-bit data model. To resume debugging, execute the ::release -a dcmd, quit MDB, and then execute **mdb** -p pid to re-attach the 64-bit debugger to the process.

If a 64-bit victim process execs a 32-bit program, then "stop" will return to the command prompt, but the debugger will only provide limited capabilities for examining the new process. All built-in dcmds will work as advertised, but loadable dcmds will not since they do not perform data model conversion of structures. The user should release and reattach the debugger to the process as described above in order to restore full debugging capabilities.

Interaction with Job Control

If the debugger is attached to a user process that is stopped by job control (that is, it stopped in response to SIGTSTP, SIGTTIN, or SIGTTOU), the process may not be able to be set running again when it is continued by a continue dcmd. If the victim process is a member of the same session (that is, it shares the same controlling terminal as MDB), MDB will attempt to bring the associated process group to the foreground and continue the process with SIGCONT to resume it from job control stop. When MDB is detached from such a process, it will restore the process group to the background before exiting. If the victim process is not a member of the same session, MDB cannot safely bring the process group to the foreground, so it will continue the process with respect to the debugger but the process will remain stopped by job control. MDB will print a warning in this case, and the user must issue a **fg** command from the appropriate shell in order to resume the process.

Process Attach and Release

When MDB attaches to a running user process, the process is stopped and remains stopped until one of the continue dcmds is applied, or the debugger quits. If the -o nostop option is enabled prior to attaching the debugger to a process with -p or prior to issuing an ::attach or :A command, MDB will attach to the process but not stop it. While the process is still running, it may be inspected as usual (albeit with inconsistent results) and breakpoints or other tracing flags may be enabled. If the :c or ::cont dcmds are executed while the process is running, the debugger will wait for the process to stop. If no traced software events occur, the user can send an interrupt (^C) after :c or ::cont to force the process to stop and return control to the debugger.

MDB releases the current running process (if any) when the :R, ::release, :r, ::run, \$q, or ::quit dcmds are executed, or when the debugger terminates as the result of an EOF or signal. If the process was originally created by the debugger using :r or ::run, it will be forcibly terminated as if by SIGKILL when it is released. If the process was already running prior to attaching MDB to it, it will be set running again when it is released. A process may be released and left stopped and abandoned using the ::release -a option.

Chapter 7. Kernel Execution Control

This chapter describes the MDB features for execution control of the live operating system kernel available when running **kmdb**. **kmdb** is a version of MDB specifically designed for kernel execution control and live kernel debugging. Using **kmdb**, the kernel can be controlled and observed in much the same way that a user process can be controlled and observed using **mdb**. The kernel execution control functionality includes instruction-level control of kernel threads executing on each CPU, enabling developers to single-step the kernel and inspect data structures in real time.

Both **mdb** and **kmdb** share the same user interface. All of the execution control functionality described in Chapter 6, Execution Control is available in **kmdb**, and is identical to the set of commands used to control user processes. The commands used to inspect kernel state, described in Chapter 3, Language Syntax and Chapter 5, Built-in Commands, are also available when using **kmdb**. Finally, the commands specific to the Solaris kernel implementation, described in Chapter 8, Kernel Debugging Modules, are available unless otherwise noted. This chapter describes the remaining features that are specific to **kmdb**.

Booting, Loading, and Unloading

To facilitate the debugging of kernel startup, **kmdb** can be loaded during the earliest stages of the boot process, before control has passed from the kernel runtime linker (krtld) to the kernel. **kmdb** may be loaded at boot using the k boot flag, the kmdb boot file, or the kadb boot file (for compatibility). If **kmdb** is loaded at boot, the debugger cannot be unloaded until the system subsequently reboots. Some functionality will not be immediately available during the earliest stages of boot. In particular, debugging modules will not be loaded until the kernel module subsystem has initialized. Processor-specific functionality will not be enabled until the kernel has completed the processor identification process.

If you boot your system using the k option, **kmdb** will automatically load during the boot process. You can use the d boot option to request a debugger breakpoint prior to starting the kernel. This feature works with the default kernel as well as alternate kernels. For example, to boot a SPARC system with **kmdb** and request immediate entry to the debugger, type any of the following commands:

```
ok boot -kd
ok boot kmdb -d
ok boot kadb -d
```

To boot an x86 system in the same manner, type any of the following commands:

```
Select (b)oot or (i)nterpreter: b -kd
Select (b)oot or (i)nterpreter: b kmdb -d
Select (b)oot or (i)nterpreter: b kadb -d
```

To boot a SPARC system with **kmdb** and load an alternate 64-bit kernel, type the following command:

```
ok boot kernel.test/sparcv9/unix -k
```

To boot an x86 system with **kmdb** and load an alternate 64-bit kernel, type the following command:

```
Select (b)oot or (i)nterpreter: b kernel.test/amd64/unix -k
```

If the boot file is set to the string kmdb or kadb and you want to boot an alternate kernel, use the D option to specify the name of the kernel to boot. To boot a SPARC system in this manner, type the following command:

```
ok boot kmdb -D kernel.test/sparcv9/unix
```

To boot a 32-bit x86 system in this manner, type the following command:

```
Select (b) or (i)nterpreter: b kmdb -D kernel.test/unix
```

To boot a 64-bit x86 system in this manner, type the following command:

```
Select (b) or (i)nterpreter: b kmdb -D kernel.test/amd64/unix
```

To debug a system that has already booted, use the **mdb** K option to load **kmdb** and stop kernel execution. When the debugger is loaded using this method, it can be subsequently unloaded. You can unload **kmdb** when you are done debugging by specifying the u option to the ::quit dcmd. Alternatively, you can resume execution of the operating system using the command mdb -U.

Terminal Handling

kmdb always uses the system console for interaction. **kmdb** will determine the appropriate terminal type according to the following rules:

- If the system being debugged uses an attached keyboard and monitor for its console and the debugger is loaded at boot, the terminal type will be determined automatically based upon the platform architecture and console terminal settings.
- If the system begin debugged uses a serial console and the debugger is loaded at boot, a default terminal type of vt100 will be assumed.
- If the debugger is loaded by running **mdb** -**K** on the console, the value of the \$TERM environment variable will be used as the terminal type.
- If the debugger is loaded by running **mdb** -**K** on a terminal that is not the console, the debugger will use the terminal type that has been configured for use with the system console login prompt.

You can use the ::term dcmd from within **kmdb** to display the terminal type.

Debugger Entry

The operating system kernel will implicitly stop executing and enter **kmdb** when a breakpoint is reached or according to the other execution control settings described in Chapter 6, Execution Control. You can use the **mdb** K option or an appropriate keyboard break sequence to request explicit entry to **kmdb**. On a SPARC system console, use the STOP-A key sequence to send a break and enter **kmdb**. On an x86 system console, use the F1–A key sequence to send a break and enter **kmdb**. You can use the **kbd** command to customize the escape sequence on your Solaris system. To enter **kmdb** on a system with a serial console, use the appropriate serial console command to send a break sequence.

Processor-Specific Features

Some **kmdb** functionality is specific to an individual processor architecture. For example, various x86 processors support a hardware branch tracing capability that is not found on some other processor architectures. Access to processor-specific features is provided through processor-specific dcmds that are only present on systems that support them. The availability of processor-specific support will be indicated in the output of the **::status** dcmd. The debugger relies upon the kernel to determine the processor type. Therefore, even though the debugger may provide features for a given processor architecture, this support will not be exposed until the kernel has progressed to the point where processor identification has completed.

Chapter 8. Kernel Debugging Modules

This chapter describes the debugger modules, dcmds, and walkers provided to debug the Solaris kernel. Each kernel debugger module is named after the corresponding Solaris kernel module, so that it will be loaded automatically by MDB. The facilities described here reflect the current Solaris kernel implementation and are subject to change in the future; writing shell scripts that depend on the output of these commands is not recommended. In general, the kernel debugging facilities described in this chapter are meaningful only in the context of the corresponding kernel subsystem implementation. See Related Books and Papers for a list of references that provide more information about the Solaris kernel implementation.

Note

MDB exposes kernel implementation details that are subject to change at any time. This guide reflects the Solaris kernel implementation as of the date of publication of this guide. Information provided in this guide about modules, dcmds, walkers, and their output formats and arguments might not be correct or applicable to past or future Solaris releases.

Generic Kernel Debugging Support (genunix) Kernel Memory Allocator

This section discusses the dcmds and walkers used to debug problems identified by the Solaris kernel memory allocator and to examine memory and memory usage. The dcmds and walkers described here are discussed in more detail in Chapter 9, Debugging With the Kernel Memory Allocator.

dcmds

thre address of a kernel thread, print a list of memory allocations it has performed in reverse **::allocations** and order.

bufce Print a summary of the bufct1 information for the specified bufct1 address. If one or more ::bufcetptions are present, the bufctl information is printed only if it matches the criteria defined by the [a address drums drums drums] running this way, the dcmd can be used as a filter for input from a pipeline. The a [c cadrition] indicates that the bufctl's corresponding buffer address must equal the specified address. The [e ear lopeist] indicates that a program counter value from the specified caller must be present in the [1 labelist] saved stack trace. The e option indicates that the bufctl's timestamp must be greater than or [t the quality of the specified earliest timestamp. The 1 option indicates that the bufctl's timestamp must be less than or equal to the specified latest timestamp. The t option indicates that the bufctl's thread pointer must be equal to the specified thread address.

[addThes:stfindleaks dcmd provides powerful and efficient detection of memory leaks in kernel crash] ::findleaks where the full set of kmem debug features has been enabled. The first execution of ::findleaks [v] processes the dump for memory leaks (this can take a few minutes), then coalesces the leaks by the allocation stack trace. The findleaks report shows a bufctl address and the topmost stack frame for each memory leak that was identified.

If the v option is specified, the dcmd prints more verbose messages as it executes. If an explicit address is specified prior to the dcmd, the report is filtered and only leaks whose allocation stack traces contain the specified function address are displayed.

threwiden the address of a kernel thread, print a list of memory frees it has performed, in reverse :: freedbyonological order.

valusearch the kernel address space for pointer-aligned addresses that contain the specified pointer-sized
::kgreplue. The list of addresses that contain matching values is then printed. Unlike MDB's built-in search operators, ::kgrep searches every segment of the kernel's address space and searches across discontiguous segment boundaries. On large kernels, ::kgrep can take a considerable amount of time to execute.

::kmaStaplay the list of kernel memory allocator caches and virtual memory arenas, along with corresponding statistics.

::kmahsersinformation about the medium and large users of the kernel memory allocator that have current [ef] memory allocations. The output consists of one entry for each unique stack trace specifying the [cachetal amount of memory and number of allocations that was made with that stack trace. This dcmd ...] requires that the KMF_AUDIT flag is set in kmem_flags.

If one or more cache names (for example, kmem_alloc_256) are specified, the scan of memory usage is restricted to those caches. By default all caches are included. If the e option is used, the small users of the allocator are included. The small users are allocations that total less than 1024 bytes of memory or for which there are less than 10 allocations with the same stack trace. If the f option is used, the stack traces are printed for each individual allocation.

[addFesssat and display the kmem_cache structure stored at the specified address, or the complete set] ::kmeinactiaellemem cache structures.

::kmeDisplogy the complete set of kmem transaction logs, sorted in reverse chronological order. This dcmd uses a more concise tabular output format than ::kmalog.

[add Wesi by the integrity of the kmem_cache structure stored at the specified address, or the complete set] ::kmeinc tierifymem_cache structures. If an explicit cache address is specified, the dcmd displays more verbose information regarding errors; otherwise, a summary report is displayed. The ::kmem_verify dcmd is discussed in more detail in Kernel Memory Caches.

[addFesssat and display the vmem structure stored at the specified address, or the complete set of active] ::vmemem structures. This structure is defined in <sys/vmem_impl.h>.

addr Examat and display the vmem_seg structure stored at the specified address. This structure is defined :: vmein_seg s / vmem_impl.h>.

addr Report information about the specified address. In particular, ::whatis will attempt to determine if ::whatise address is a pointer to a kmem-managed buffer or another type of special memory region, such as [abv]a thread stack, and report its findings. If the a option is present, the dcmd reports all matches instead of just the first match to its queries. If the b option is present, the dcmd also attempts to determine if the address is referred to by a known kmem bufctl. If the v option is present, the dcmd reports its progress as it searches various kernel data structures.

Walkers

allocdby Given the address of a kthread_t structure as a starting point, iterate over the set

of bufctl structures corresponding to memory allocations performed by this kernel

thread.

bufctl Given the address of a kmem_cache_t structure as a starting point, iterate over the

set of allocated bufctls associated with this cache.

freectl Given the address of a kmem_cache_t structure as a starting point, iterate over the

set of free bufctls associated with this cache.

freedby Given the address of a kthread_t structure as a starting point, iterate over the set of

bufctl structures corresponding to memory deallocations performed by this kernel

thread.

freemem Given the address of a kmem_cache_t structure as a starting point, iterate over the

set of free buffers associated with this cache.

kmem Given the address of a kmem_cache_t structure as a starting point, iterate over the

set of allocated buffers associated with this cache.

<sys/kmem_impl.h>.

kmem_cpu_cache Given the address of a kmem_cache_t structure as a starting point, iterate over the

per-CPU kmem_cpu_cache_t structures associated with this cache. This structure

is defined in <sys/kmem_impl.h>.

kmem_slab Given the address of a kmem_cache_t structure as a starting point, iterate over

the set of associated kmem_slab_t structures. This structure is defined in <sys/

kmem_impl.h>.

leak Given the address of a bufctl structure, iterate over the set of bufctl structures

corresponding to leaked memory buffers with similar allocation stack traces. The ::findleaks dcmd must be applied to locate memory leaks before the leak walker

can be used

leakbuf Given the address of a bufctl structure, iterate over the set of buffer addresses

corresponding to leaked memory buffers with similar allocation stack traces. The ::findleaks dcmd must be applied to locate memory leaks before the leakbuf

walker can be used.

File Systems

The MDB file systems debugging support includes a built-in facility to convert vnode pointers to the corresponding file system path name. This conversion is performed using the Directory Name Lookup Cache (DNLC); because the cache does not hold all active vnodes, some vnodes might not be able to be converted to path names and "??" is displayed instead of a name.

dcmds

::fsinfo Display a table of mounted file systems, including the vfs_t

address, ops vector, and mount point of each file system.

::Iminfo Display a table of vnodes with active network locks registered with

the lock manager. The pathname corresponding to each vnode is

shown.

address::vnode2path [v] Display the pathname corresponding to the given vnode address. If

the v option is specified, the dcmd prints a more verbose display, including the vnode pointer of each intermediate path component.

Walkers

buf Iterate over the set of active block I/O transfer structures (buf_t structures). The buf structure is defined in <sys/buf.h> and is described in more detail in buf(9S).

Virtual Memory

This section describes the debugging support for the kernel virtual memory subsystem.

dcmds

address ::addr2smap Print the smap structure address that corresponds to the given

[offset] address in the kernel's segmap address space segment.

as ::as2proc Display the proc_t address for the process corresponding to the

as_t address as.

[address] ::memlist [aiv] Display the specified memlist structure or one of the well-known

memlist structures. If no memlist address and options are present or if the i option is present, the memlist representing physically installed memory is displayed. If the a option is present, the memlist representing available physical memory is displayed. If the v option is present, the memlist representing available virtual memory is

displayed.

::memstat Display a system-wide memory usage summary. The amount and

percentage of system memory consumed by different classes of pages (kernel, anonymous memory, executables and libraries, page cache, and free lists) are displayed, along with the total amount of

system memory.

[address] ::page Display the properties of the specified page_t. If no page_t

address is specified, the dcmd displays the properties of all system

pages.

seg::seg Format and display the specified address space segment (seg_t

address).

[address] ::swapinfo Display information on all active swapinfo structures or about the

specified struct swapinfo. The vnode, filename, and statistics for

each structure are displayed.

vnode::vnode2smap [offset] Print the smap structure address that corresponds to the given

vnode_t address and offset.

Walkers

anon Given the address of an anon_map structure as a starting point, iterate over the set of related

anon structures. The anon map implementation is defined in <vm/anon.h>.

conjunction with the ::memlist dcmd to display each span.

page Iterate over all system page structures. If an explicit address is specified for the walk, this is

taken to be the address of a vnode and the walker iterates over only those pages associated

with the vnode.

seg Given the address of an as_t structure as a starting point, iterate over the set of address space

segments (seg structures) associated with the specified address space. The seg structure is

defined in <vm/seg.h>.

swapinfo Iterate over the list of active swapinfo structures. This walker may be used in conjunction

with the ::swapinfo dcmd.

CPUs and the Dispatcher

This section describes the facilities for examining the state of the cpu structures and the kernel dispatcher.

dcmds

::callout Display the callout table. The function, argument, and expiration time

for each callout is displayed.

::class Display the scheduling class table.

[cpuid]::cpuinfo[v] Display a table of the threads currently executing on each CPU. If

an optional CPU ID number or CPU structure address is specified prior to the dcmd name, only the information for the specified CPU is displayed. If the v option is present, ::cpuinfo also displays the runnable threads waiting to execute on each CPU as well as the active

interrupt threads.

Walkers

cpu Iterate over the set of kernel CPU structures. The cpu t structure is defined in <sys/cpuvar.h>.

Device Drivers and DDI Framework

This section describes dcmds and walkers that are useful for kernel developers as well as third-party device driver developers.

dcmds

addressen the address of a kernel name-to-major number binding hash table entry (struct bind), display ::bindingnassbirding name, major number, and pointer to the next element.

::devblidgings he list of all instances of the named driver. The output consists of an entry for each instance, devibeginning with the pointer to the struct dev_info (viewable with \$<devinfo or ::devinfo), the namedriver name, the instance number, and the driver and system properties associated with that instance.

addr Prist the system and driver properties associated with a devinfo node. If the q option is specified, ::devinfoy a quick summary of the device node is shown.

```
]
[ d
```

addr Print the name of the driver (if any) associated with the devinfo node.

::devinfo2driver

[addDissplay the kernel's devnames table along with the dn_head pointer, which points at the driver] ::devnstance list. If the v flag is specified, additional information stored at each entry in the devnames [v table is displayed.

[dev Display the kernel device tree starting at the device node specified by devinfo. If devinfo is not] ::prtprovided, the root of the device tree is assumed by default. If the c option is specified, only children [cpvof the given device node are displayed. If the p option is specified, only ancestors of the given device node are displayed. If v is specified, the properties associated with each node are displayed.

[majDisplay the driver name corresponding to the specified major number. The major number can be num specified as an expression preceding the dcmd or as a command-line argument.

]::major2name

```
[major-
num
]
```

[addPristsall of the device nodes that correspond to the specified model address.

]::modctl2devinfo

::nam@2majorlevice driver name, display its major number.

```
driver-
name
```

[add Sieves a softstate state pointer (see ddi_soft_state_init(9F)) and a device instance number, display] ::softstate for that instance.

```
[instance-
number
```

Walkers

binding_hash Given the address of an array of kernel binding hash table entries (struct bind

**), walk all entries in the hash table and return the address of each struct bind.

devinfo First, iterate over the parents of the given devinfo and return them in order

of seniority from most to least senior. Second, return the given devinfo itself. Third, iterate over the children of the given devinfo in order of seniority from most to least senior. The dev_info struct is defined in <sys/

ddi_impldefs.h>.

devinfo_children First, return the given devinfo, then iterate over the children of the given devinfo

in order of seniority from most to least senior. The dev_info struct is defined

in <sys/ddi_impldefs.h>.

least senior, and then return the given devinfo. The dev_info struct is defined

in <sys/ddi_impldefs.h>.

devi_next Iterate over the siblings of the given devinfo. The dev_info struct is defined in

<sys/ddi_impldefs.h>.

devnames Iterate over the entries in the devnames array. This structure is defined in

<sys/autoconf.h>.

softstate Given a softstate pointer (see **ddi_soft_state_init**(9F)) display all non-NULL

pointers to driver state structures.

softstate_all Given a softstate pointer (see **ddi_soft_state_init**(9F)) display all pointers to

driver state structures. Note that the pointers for unused instances will be NULL.

STREAMS

This section describes dcmds and walkers that are useful for kernel developers as well as developers of third-party STREAMS modules and drivers.

dcmds

address of the address of an mblk_t, print the address of the corresponding dblk_t.

::mblk2dblk

[add:Nexisty the integrity of one or more message blocks. If an explicit message block address is specified, ::mblkheverifyrity of this message block is checked. If no address is specified, the integrity of all active message blocks are checked. This dcmd produces output for any invalid message block state that is detected.

addr Sissen the address of a queue_t, print the address of the corresponding syncq_t data structure.

::q2syncq

address of the peer read or write queue structure.

::q2otherq

addressen the address of a queue_t, print the address of the corresponding read queue.

::q2rdq

address of the address of a queue t, print the address of the corresponding write queue.

::q2wrq

[addDissplay a visual picture of a kernel STREAM data structure, given the address of the stdata_t] ::streameture representing the STREAM head. The read and write queue pointers, byte count, and flags for each module are shown, and in some cases additional information for the specific queue is shown in the margin.

addr Eister and display the specified syncq_t data structure. With no options, various properties of the ::syncqncq_t are shown. If the v option is present, the syncq flags are decoded in greater detail. If the f, [v] F, t, or T options are present, the syncq is displayed only if it matches the criteria defined by the [f flaggliments to these options; in this way, the dcmd can be used as a filter for input from a pipeline. The

[F flag]ption indicates that the specified flag (one of the SQ_ flag names from <sys/strsubr.h>)
[t typhe]st be present in the syncq flags. The F option indicates that the specified flag must be absent from
[T typhe]syncq flags. The t option indicates that the specified type (one of the SQ_CI or SQ_CO type names from <sys/strsubr.h>) must be present in the syncq type bits. The T option indicates that the specified type must be absent from the syncq type bits.

addressen the address of a syncq_t, print the address of the corresponding queue_t data structure. ::syncq2q

Walkers

b_cont	Given the address of an mblk_t, iterate over the set of associated message structures by following the b_cont pointer. The b_cont pointer is used to link a given message block to the next associated message block that is the continuation of the same message. The message block is described in more detail in msgb(9S)
b_next	Given the address of an mblk_t, iterate over the set of associated message structures by following the b_next pointer. The b_next pointer is used to link a given message block to the next associated message block on a given queue. The message block is described in more detail in msgb(9S).
qlink	Given the address of a queue_t structure, walk the list of related queues using the q_link pointer. This structure is defined in <sys stream.h="">.</sys>
qnext	Given the address of a queue_t structure, walk the list of related queues using the q_next pointer. This structure is defined in <sys stream.h="">.</sys>
readq	Given the address of an stdata_t structure, walk the list of read-side queue structures.
writeq	Given the address of an stdata_t structure, walk the list of write-side queue structures.

Circum the address of an arbitrate area the act of associated asso

Networking

The following dcmds and walkers are provided to help debug the core kernel networking stack protocols.

dcmds

addr Sissen the address of a kernel MI_O, filter and display the MI_O or its payload. If the p option is
 ::mi specified, then the address of the corresponding payload of the MI_O is displayed, otherwise the
 [p] MI_O itself is displayed. Specifying filter d or m enables the dcmd to filter device or module MI_O
 [d | objects respectively.
 m]

::netsfiltow network statistics and active connections. If the a option is present, the state of all sockets is [av] displayed. If the v option is present, more verbose output is displayed. If the f option is present, only [finetenections associated with the specified address family are displayed. If the p option is present, | inebfly connections associated with the specified protocols are displayed.

unix]
[Ptcp
udp]

[addFitses and displays sonode objects. If no address is given, then the list of AF_UNIX sockets is]::sontide ayed, otherwise only the specified sonode is displayed. If the f option is present, then only [finsackets of the given family will be output. If the t option is present, then only sonodes of the

| inegiven type will be output. If the p option is present, then only sockets of the given protocol will | unite displayed.

|id]
[t stream
|dgram
|raw
|id]

dcmd will be verbose.

[p id]

[addFitses and displays tcpb objects. If no address is specified, all connections are walked, otherwise] ::tcphnly the specified tcpb is filtered/displayed. Specifying a filters for only active connections and P can [av] be used to filter for TCP IPv4 or IPv6 connections. The tcpb dcmd is intelligent about filtering TCP connections, and if a IPv6 TCP connection is in a state that would still facilitate a IPv4 connection, v4 the P filter considers the connection as both IPv4 and IPv6 in much the same way that ::netstat | v6] does. If the dcmd is not being used as a filter and the v option is specified, then the output of the

Walkers

ar	Given the address of an ar, this walker walks all ar objects from the given ar to the final ar. If no address is specified, all ar objects are walked.
icmp	Given the address of an icmp, this walker walks all icmp objects from the given icmp to the final icmp. If no address is specified, all icmp objects are walked.
ill	Given the address of an interface link layer structure (ill), this walker walks all ill objects from the given ill to the final. If no address is specified, all ill objects are walked.
ipc	Given the address of an ipc, this walker walks all ipc objects from the given ipc to the final ipc. If no address is specified, all ipc objects are walked.
mi	Given the address of a MI_O, walk all the MI_O's in this MI.
sonode	Given the address of a AF_UNIX sonode, walk the associated list of AF_UNIX sonodes beginning with the given sonode. If no address is specified, this walker walks the list of all AF_UNIX sockets.
tcpb	beginning with the given sonode. If no address is specified, this walker walks the list of all

Files, Processes, and Threads

This section describes dcmds and walkers used to format and examine various fundamental file, process, and thread structures in the Solaris kernel.

dcmds

process the file_t address corresponding to the file descriptor fd-num associated with the specified
::fd process. The process is specified using the virtual address of its proc_t structure.
fd-

num

thre Print the stack trace associated with the given kernel thread, identified by the virtual address of its ::findkthckad_t structure. The dcmd employs several different algorithms to locate the appropriate stack [combackdrace. If an optional command string is specified, the dot variable is reset to the frame pointer address of the topmost stack frame, and the specified command is evaluated as if it had been typed at the command line. The default command string is "<.\$C0"; that is, print a stack trace including frame pointers but no arguments.</p>

::pgrepisplay process information for processes whose name matches the regexp regular expression [x] [nimitern. The ::pgrep dcmd is similar to the pgrep(1) command. The ::pgrep dcmd is used to pattern regematch against all processes. When the n option is used, display only the newest process that matches the pattern. When the o option is used, display only the oldest process that matches the pattern. When the x option is used, display only those processes whose names are exactly the same as the search pattern.

In kmdb(1), the regexp used with ::pgrep must be a plain alpha-numeric text string.

pid Print the proc_t address corresponding to the specified PID. Recall that MDB's default base is ::pid2pexxdecimal, so decimal PIDs obtained using pgrep(1) or ps(1) should be prefixed with 0t.

processing the memory map of the process indicated by the given process address. The dcmd displays ::pmaputput using a format similar to pmap(1). If the q option is present, the dcmd displays an abbreviated [q] form of its output that requires less processing time.

[addPrists summary of the information related to the specified process, or all active system processes,] ::ps similar to ps(1). If the f option is specified, the full command name and initial arguments are printed. [fltTp]he l option is specified, the LWPs associated with each process are printed. If the t option is specified, the kernel threads associated with each process LWP are printed. If the T option is specified, the task ID associated with each process is displayed. If the P option is specified, the project ID associated with each process is displayed.

::ptrePrint a process tree, with child processes indented from their respective parent processes. The dcmd displays output using a format similar to ptree(1).

address a list of the active kernel task structures and their associated ID numbers and attributes. The ::taskprocess task ID is described in more detail in settaskid(2).

[addDissplay properties of the specified kernel kthread_t structure. If no kthread_t address is specified,] ::thrteadproperties of all kernel threads are displayed. The dcmd options are used to control which output [bdfixplass]ns are shown. If no options are present, the i option is enabled by default. If the b option is present, information relating to the thread's turnstile and blocking synchronization object is shown. If the d option is present, the thread's dispatcher priority, binding, and last dispatch time is shown. If the f option is present, threads whose state is TS_FREE are elided from the output. If the i option is present (the default), thread state, flags, priority, and interrupt information is shown. If the m option is present, all of the other output options are merged together on to a single output line. If the p option is present, the thread's process, LWP, and credential pointers are displayed. If the s option is present, the thread's signal queue and masks of pending and held signals are shown.

vnod Given a vnode_t address, print the proc_t addresses of all processes that have this vnode currently ::whereperin their file table.

Walkers

file Given the address of a proc_t structure as a starting point, iterate over the set of open files (file_t structures) associated with the specified process. The file_t structure is defined in <sys/file.h>.

proc Iterate over the active process (proc_t) structures. This structure is defined in <sys/

proc.h>.

task Given a task pointer, iterate over the list of proc_t structures for processes that are members

of the given task.

thread Iterate over a set of kernel thread (kthread_t) structures. If the global walk is invoked, all kernel

threads are returned by the walker. If a local walk is invoked using a proc_t address as the starting point, the set of threads associated with the specified process is returned. The kthread_t

structure is defined in <sys/thread.h>.

Synchronization Primitives

This section describes dcmds and walkers used to examine particular kernel synchronization primitives. The semantics of each primitive are discussed in the corresponding (9f) section of the manual pages.

dcmds

rwlock::rwlock Given the address of a readers-writers lock (see rwlock(9F)),

display the current state of the lock and the list of waiting threads.

address::sobj2ts Convert the address of a synchronization object to the address of

the corresponding turnstile and print the turnstile address.

[address] ::turnstile Display the properties of the specified turnstile_t. If no turnstile_t

address is specified, the dcmd displays the properties of all

turnstiles.

[address] ::wchaninfo [v] Given the address of a condition variable (see condvar(9F)) or

semaphore (see semaphore(9F)), display the current number of waiters on this object. If no explicit address is specified, display all such objects that have waiting threads. If the v option is specified,

display the list of threads that are blocked on each object.

Walkers

blocked Given the address of a synchronization object (such as a mutex(9F) or rwlock(9F)), iterate

over the list of blocked kernel threads.

wchan Given the address of a condition variable (see condvar(9F)) or semaphore (see

semaphore(9F)), iterate over the list of blocked kernel threads.

Cyclics

The cyclic subsystem is a low-level kernel subsystem that provides high resolution, per-CPU interval timer facilities to other kernel services and programming interfaces.

dcmds

::cycinfo [vV] Display the cyclic subsystem per-CPU state for each CPU. If the v option is

present, a more verbose display is shown. If the V option is present, an even

more verbose display than v is shown.

address ::cyclic Format and display the cyclic_t at the specified address.

::cyccover Display cyclic subsystem code coverage information. This information is

available only in a DEBUG kernel.

::cyctrace Display cyclic subsystem trace information. This information is available only

in a DEBUG kernel.

Walkers

cyccpu Iterate over the per-CPU cyc_cpu_t structures. This structure is defined in <sys/

cyclic_impl.h>.

cyctrace Iterate over the cyclic trace buffer structures. This information is only available in a DEBUG

kernel.

Task Queues

The task queue subsystem provides general-purpose asynchronous task scheduling for a variety of clients in the kernel.

dcmds

address ::taskq_entry Print the contents of the specified struct taskq_entry.

Walkers

taskq_entry Given the addresss of a taskq structure, iterate over the list of taskq_entry structures.

Error Queues

The error queue subsystem provides general-purpose asynchronous error event processing for platform-specific error handling code.

dcmds

[address] ::errorq Display a summary of information relating to the specified error queue.

If no address is given, display information relating to all system error queues. The address, name, queue length, and data element size for each

queue are displayed, along with various queue statistics.

Walkers

errorg Walk the list of system error queues and return the address of each individual error

queue.

errorq_data Given the address of an error queue, return the address of each pending error event

data buffer.

Configuration

This section describes demds that can be used to examine system configuration data.

dcmds

::system

Display the contents of the system(4) configuration file at the time the kernel parsed the file during system initialization.

Interprocess Communication Debugging Support (ipc)

The ipc module provides debugging support for the implementation of the message queue, semaphore, and shared memory interprocess communication primitives.

dcmds

::ipcs [1] Display a listing of system-wide IPC identifiers, corresponding

to known message queues, semaphores, and shared memory segments. If the 1 option is specified, a longer listing of information

is shown.

address::msg[1][t type] Display the properties of the specified message queue element

(struct msg). If the 1 option is present, the raw contents of the message are displayed in hexadecimal and ASCII. If the t option is present, it can be used to filter the output and only display messages of the specified type. This can be useful when piping the output of

the msgqueue walker to ::msg.

id::msqid [k] Convert the specified message queue IPC identifier to a pointer

to the corresponding kernel implementation structure and print the address of this kernel structure. If the k option is present, the id is instead interpreted as a message queue key to match (see

msgget(2)).

[address] ::msqid_ds [1] Print the specified msqid_ds structure or a table of the active

msqid_ds structures (message queue identifiers). If the 1 option is

specified, a longer listing of information is displayed.

id::semid [k] Convert the specified semaphore IPC identifier to a pointer to

the corresponding kernel implementation structure and print the address of this kernel structure. If the k option is present, the *id* is instead interpreted as a semaphore key to match (see semget(2)).

instead interpreted as a semaphore key to match (see semger(2)

[address] ::semid_ds [1] Print the specified semid_ds structure or a table of the active

semid_ds structures (semaphore identifiers). If the 1 option is

specified, a longer listing of information is displayed.

id::shmid [k] Convert the specified shared memory IPC identifier to a pointer

to the corresponding kernel implementation structure and print the address of this kernel structure. If the k option is present, the id is instead interpreted as a shared memory key to match (see

shmget(2)).

[address] ::shmid_ds [1] Print the specified shmid_ds structure or a table of the active

shmid_ds structures (shared memory segment identifiers). If the 1 option is specified, a longer listing of information is displayed.

Walkers

msg Walk the active msqid_ds structures corresponding to message queue identifiers. This

structure is defined in <sys/msg.h>.

msgqueue Iterate over the message structures that are currently enqueued on the specified message

queue.

sem Walk the active semid_ds structures corresponding to semaphore identifiers. This structure

is defined in <sys/sem.h>.

shm Walk the active shmid_ds structures corresponding to shared memory segment identifiers.

This structure is defined in <sys/shm.h>.

Loopback File System Debugging Support (lofs)

The lofs module provides debugging support for the lofs(7FS) file system.

dcmds

[address] ::Inode Print the specified lnode_t, or a table of the active lnode_t structures in

the kernel.

address::lnode2dev Print the dev_t (vfs_dev) for the underlying loopback mounted filesystem

corresponding to the given lnode_t address.

address::Inode2rdev Print the dev_t (li_rdev) for the underlying loopback mounted file system

corresponding to the given lnode t address.

Walkers

lnode Walk the active lnode_t structures in the kernel. This structure is defined in <sys/fs/
lofs node.h>.

Internet Protocol Module Debugging Support (ip)

The ip module provides debugging support for the ip(7P) driver

dcmds

[address] :: ire [q] Print the specified ire_t, or a table of the active ire_t structures in the

kernel. If the q flag is specified, the send and receive queue pointers are

printed instead of the source and destination addresses.

Walkers

ire Walk the active ire (Internet Route Entry) structures in the kernel. This structure is defined in <inet/ip.h>.

Kernel Runtime Link Editor Debugging Support (krtld)

This section describes the debugging support for the kernel runtime link editor, which is responsible for loading kernel modules and drivers.

dcmds

[address] ::modctl Print the specified modctl, or a table of the active modctl structures in

the kernel.

address::modhdrs Given the address of a modetl structure, print the module's ELF

executable header and section headers.

::modinfo Print information about the active kernel modules, similar to the output

of the /usr/sbin/modinfo command.

Walkers

modetl Walk the list of active modetl structures in the kernel. This structure is defined in <sys/

modctl.h>.

USB Framework Debugging Support (uhci)

The uchi module provides debugging support for the host controller interface portion of the Universal Serial Bus (USB) framework.

dcmds

address::uhci_qh [bd] Given the address of a USB UHCI controller Queue Head (QH)

structure, print the contents of the structure. If the b option is present iterate over the link_ptr chain, printing all QHs found. If the d option is present, iterate over the element_ptr chain, printing all

TDs found.

address::uhci td [d] Given the address of a USB UHCI controller Transaction Descriptor

(TD) structure, print the contents of the structure. Note this only works for Control and Interrupt TDs. If the d option is present, iterate over

the element ptr chain, printing all TDs found.

Walkers

uhci_qh Given the address of a USB UHCI controller Queue Head (QH) structure, iterate over the

list of such structures.

uhci_td Given the address of a USB UHCI controller Queue Head Descriptor (TD) structure, iterate

over the list of such structures.

USB Framework Debugging Support (usba)

The usba module provides debugging support for the platform-independent Universal Serial Bus (USB) framework.

dcmds

::usba_debug_buf Print the USB debugging information buffer.

::usba_clear_debug_buf Empty the USB debugging information buffer.

[address]::usba_device [pv] Given the address of a usba_device structure, print summary

information. If no address is supplied, this dcmd walks the global list of usba_device structures. If the p option is present, also list information for all open pipes on this device. If the v option is

present, list verbose information for each device.

address::usb_pipe_handle Given the address of a USB pipe handle structure (struct

usba_ph_impl), print summary information for this handle.

Walkers

usba_list_entry Given the address of a usba_list_entry structure, iterate over the chain of such

structures.

usba_device Walk the global list of usba_device_t structures.

usb_pipe_handle Given a usba_device_t address, walk USB pipe handles.

x86 Platform Debugging Support (unix)

These dcmds and walkers are specific to x86 platforms.

dcmds

[cpuDisplay trap trace records in reverse chronological order. The trap trace facility is available only in | addDeBUG kernels. If an explicit dot value is specified, this is interpreted as either a CPU ID number]::ttrace trap trace record address, depending on the precise value. If a CPU ID is specified, the output is [x] restricted to the buffer from that CPU. If a record address is specified, only that record is formatted. If the x option is specified, the complete raw record is displayed.

Walkers

ttrace Walk the list of trap trace record addresses in reverse chronological order. The trap trace facility

is available only in DEBUG kernels.

sun4u Platform Debugging Support (unix)

These dcmds and walkers are specific to the SPARC sun4u platform.

dcmds

[address] ::softint Display the soft interrupt vector structure at the specified address, or

display all the active soft interrupt vectors. The pending count, PIL, argument, and handler function for each structure is displayed.

::ttctl Display trap trace control records. The trap trace facility is available

only in DEBUG kernels.

[cpuid]::ttrace [x] Display trap trace records in reverse chronological order. The trap trace

facility is available only in DEBUG kernels. If an explicit dot value is specified, this is interpreted as a CPU ID number, and the output is restricted to the buffer from that CPU. If the x option is specified, the

complete raw record is displayed.

[address] ::xc_mbox Display the cross-call mailbox at the specified address, or format all

the cross-call mailboxes that have pending requests.

::xctrace Format and display cross-call trace records in reverse chronological

order that are related to CPU cross-call activity. The cross-call trace

facility is available only in DEBUG kernels.

Walkers

softint Iterate over the soft interrupt vector table entries.

ttrace Iterate over the trap trace record addresses in reverse chronological order. The trap trace

facility is only available in DEBUG kernels.

Chapter 9. Debugging With the Kernel Memory Allocator

The Solaris kernel memory (kmem) allocator provides a powerful set of debugging features that can facilitate analysis of a kernel crash dump. This chapter discusses these debugging features, and the MDB dcmds and walkers designed specifically for the allocator. Bonwick (see Related Books and Papers) provides an overview of the principles of the allocator itself. Refer to the header file <sys/kmem_impl.h> for the definitions of allocator data structures. The kmem debugging features can be enabled on a production system to enhance problem analysis, or on development systems to aid in debugging kernel software and device drivers.

Note

MDB exposes kernel implementation details that are subject to change at any time. This guide reflects the Solaris kernel implementation as of the date of publication of this guide. Information provided in this guide about the kernel memory allocator might not be correct or applicable to past or future Solaris releases.

Getting Started: Creating a Sample Crash Dump

This section shows you how to obtain a sample crash dump, and how to invoke MDB in order to examine it.

Setting kmem_flags

The kernel memory allocator contains many advanced debugging features, but these are not enabled by default because they can cause performance degradation. In order to follow the examples in this guide, you should turn on these features. You should enable these features only on a test system, as they can cause performance degradation or expose latent problems.

The allocator's debugging functionality is controlled by the kmem_flags tunable. To get started, make sure kmem_flags is set properly:

```
# mdb -k
> kmem_flags/X
kmem_flags:
kmem_flags: f
```

If kmem_flags is not set to 'f', you should add the line:

```
set kmem_flags=0xf
```

to /etc/system and reboot the system. When the system reboots, confirm that kmem_flags is set to 'f'. Remember to remove your /etc/system modifications before returning this system to production use.

Forcing a Crash Dump

The next step is to make sure crash dumps are properly configured. First, confirm that **dumpadm** is configured to save kernel crash dumps and that **savecore** is enabled. See dumpadm(1M) for more information on crash dump parameters.

Debugging With the Kernel Memory Allocator

dumpadm

```
Dump content: kernel pages
Dump device: /dev/dsk/c0t0d0s1 (swap)
Savecore directory: /var/crash/testsystem
Savecore enabled: yes
```

Next, reboot the system using the '-d' flag to reboot(1M), which forces the kernel to panic and save a crash dump.

When the system reboots, make sure the crash dump succeeded:

```
$ cd /var/crash/testsystem
$ ls
bounds unix.0 unix.1 vmcore.0 vmcore.1
```

If the dump is missing from your dump directory, it could be that the partition is out of space. You can free up space and run savecore(1M) manually as root to subsequently save the dump. If your dump directory contains multiple crash dumps, the one you just created will be the unix. [n] and vmcore. [n] pair with the most recent modification time.

Starting MDB

Now, run **mdb** on the crash dump you created, and check its status:

```
$ mdb unix.1 vmcore.1
Loading modules: [ unix krtld genunix ip nfs ipc ]
> ::status
debugging crash dump vmcore.1 (32-bit) from testsystem
operating system: 5.10 Generic (sun4u)
panic message: forced crash dump initiated at user request
```

In the examples presented in this guide, a crash dump from a 32-bit kernel is used. All of the techniques presented here are applicable to a 64-bit kernel, and care has been taken to distinguish pointers (sized differently on 32- and 64-bit systems) from fixed-sized quantities, which are invariant with respect to the kernel data model.

An UltraSPARC workstation was used to generate the example presented. Your results can vary depending on the architecture and model of system you use.

Allocator Basics

The kernel memory allocator's job is to parcel out regions of virtual memory to other kernel subsystems (these are commonly called *clients*). This section explains the basics of the allocator's operation and introduces some terms used later in this guide.

Buffer States

The functional domain of the kernel memory allocator is the set of *buffers* of virtual memory that make up the kernel heap. These buffers are grouped together into sets of uniform size and purpose, known as *caches*. Each cache contains a set of buffers. Some of these buffers are currently *free*, which means that they have not yet been allocated to any client of the allocator. The remaining buffers are *allocated*, which means that a pointer to that buffer has been provided to a client of the allocator. If no client of the allocator holds a pointer to an allocated buffer, this buffer is said to be *leaked*, because it cannot be freed. Leaked buffers indicate incorrect code that is wasting kernel resources.

Transactions

A kmem *transaction* is a transition on a buffer between the allocated and free states. The allocator can verify that the state of a buffer is valid as part of each transaction. Additionally, the allocator has facilities for logging transactions for post-mortem examination.

Sleeping and Non-Sleeping Allocations

Unlike the Standard C Library's malloc(3C) function, the kernel memory allocator can block (or *sleep*), waiting until enough virtual memory is available to satisfy the client's request. This is controlled by the 'flag' parameter to kmem_alloc(9F). A call to kmem_alloc(9F) which has the KM_SLEEP flag set can never fail; it will block forever waiting for resources to become available.

Kernel Memory Caches

The kernel memory allocator divides the memory it manages into a set of *caches*. All allocations are supplied from these caches, which are represented by the kmem_cache_t data structure. Each cache has a fixed *buffer size*, which represents the maximum allocation size satisfied by that cache. Each cache has a string name indicating the type of data it manages.

Some kernel memory caches are special purpose and are initialized to allocate only a particular kind of data structure. An example of this is the "thread_cache," which allocates only structures of type kthread_t. Memory from these caches is allocated to clients by the kmem_cache_alloc function and freed by the kmem_cache_free function.

Note

kmem_cache_alloc and kmem_cache_free are not public DDI interfaces. Do NOT write code that relies on them, because they are subject to change or removal in future releases of Solaris.

Caches whose name begins with "kmem_alloc_" implement the kernel's general memory allocation scheme. These caches provide memory to clients of kmem_alloc(9F) and kmem_zalloc(9F). Each of these caches satisfies requests whose size is between the buffer size of that cache and the buffer size of the next smallest cache. For example, the kernel has kmem_alloc_8 and kmem_alloc_16 caches. In this case, the kmem_alloc_16 cache handles all client requests for 9-16 bytes of memory. Remember that the size of each buffer in the kmem_alloc_16 cache is 16 bytes, regardless of the size of the client request. In a 14 byte request, two bytes of the resulting buffer are unused, since the request is satisfied from the kmem_alloc_16 cache.

The last set of caches are those used internally by the kernel memory allocator for its own bookkeeping. These include those caches whose names start with "kmem_magazine_" or "kmem_va_", the kmem slab cache, the kmem bufctl cache and others.

Kernel Memory Caches

This section explains how to find and examine kernel memory caches. You can learn about the various kmem caches on the system by issuing the ::kmastat command.

> ::kmastat						
cache	buf	buf	buf	memory	alloc	alloc
name	size	in use	total	in use	succeed	fail
kmem_magazine_1	8	24	1020	8192	24	0
kmem_magazine_3	16	141	510	8192	141	0
kmem_magazine_7	32	96	255	8192	96	0
• • •						
kmem_alloc_8	8	3614	3751	90112	9834113	0
kmem_alloc_16	16	2781	3072	98304	8278603	0
kmem_alloc_24	24	517	612	24576	680537	0
kmem_alloc_32	32	398	510	24576	903214	0
kmem_alloc_40	40	482	584	32768	672089	0
•••						
thread_cache	368	107	126	49152	669881	0
lwp_cache	576	107	117	73728	182	0
turnstile_cache	36	149	292	16384	670506	0
cred_cache	96	6	73	8192	2677787	0

If you run ::kmastat you get a feel for what a "normal" system looks like. This will help you to spot excessively large caches on systems that are leaking memory. The results of ::kmastat will vary depending on the system you are running on, how many processes are running, and so forth.

Another way to list the various kmem caches is with the ::kmem_cache command:

> ::kmem_	_cache				
ADDR	NAME	FLAG	CFLAG	BUFSIZE	BUFTOTL
70036028	kmem_magazine_1	0020	0e0000	8	1020
700362a8	kmem_magazine_3	0020	0e0000	16	510
70036528	kmem_magazine_7	0020	0e0000	32	255
70039428	kmem_alloc_8	020f	000000	8	3751
700396a8	kmem_alloc_16	020f	000000	16	3072
70039928	kmem_alloc_24	020f	000000	24	612
70039ba8	kmem_alloc_32	020f	000000	32	510
7003a028	kmem_alloc_40	020f	000000	40	584

This command is useful because it maps cache names to addresses, and provides the debugging flags for each cache in the FLAG column. It is important to understand that the allocator's selection of debugging features is derived on a per-cache basis from this set of flags. These are set in conjunction with the global kmem_flags variable at cache creation time. Setting kmem_flags while the system is running has no effect on the debugging behavior, except for subsequently created caches (which is rare after boot-up).

Next, walk the list of kmem caches directly using MDB's kmem_cache walker:

```
> ::walk kmem_cache
70036028
700362a8
```

70036528 700367a8

. . .

This produces a list of pointers that correspond to each kmem cache in the kernel. To find out about a specific cache, apply the kmem_cache macro:

> 0x70039928\$ <kmem_cache< th=""></kmem_cache<>			
0x70039928:	lock		
0x70039928:	owner/waiters		
	0		
0x70039930:	flags	freelist	offset
	20f	707c86a0	24
0x7003993c:	global_alloc	global_free	alloc_fail
	523	0	0
0x70039948:	hash_shift	hash_mask	hash_table
	5	1ff	70444858
0x70039954:	nullslab		
0x70039954:	cache	base	next
	70039928	0	702d5de0
0x70039960:	prev	head	tail
	707c86a0	0	0
0x7003996c:	refcnt	chunks	
	-1	0	
0x70039974:	constructor	destructor	reclaim
	0	0	0
0x70039980:	private	arena	cflags
	0	104444f8	0
0x70039994:	bufsize	align	chunksize
	24	8	40
0x700399a0:	slabsize	color	maxcolor
	8192	24	32
0x700399ac:	slab_create	slab_destroy	buftotal
	3	0	612
0x700399b8:	bufmax	rescale	lookup_depth
	612	1	0
0x700399c4:	kstat	next	prev
	702c8608	70039ba8	700396a8
0x700399d0:	name kmem_al		
0x700399f0:	bufctl_cache	magazine_cache	magazine_size
	70037ba8	700367a8	15

Important fields for debugging include 'bufsize', 'flags' and 'name'. The name of the kmem_cache (in this case "kmem_alloc_24") indicates its purpose in the system. Bufsize indicates the size of each buffer in this cache; in this case, the cache is used for allocations of size 24 and smaller. 'flags' indicates what debugging features are turned on for this cache. You can find the debugging flags listed in <sys/kmem_impl.h>. In this case 'flags' is 0x20f, which is KMF_AUDIT | KMF_DEADBEEF | KMF_REDZONE | KMF_CONTENTS | KMF_HASH. This document explains each of the debugging features in subsequent sections.

When you are interested in looking at buffers in a particular cache, you can walk the allocated and freed buffers in that cache directly:

```
> 0x70039928::walk kmem
```

```
704ba010
702ba008
704ba038
702ba030
...
> 0x70039928::walk freemem
70a9ae50
70a9ae28
704bb730
704bb2f8
```

MDB provides a shortcut to supplying the cache address to the kmem walker: a specific walker is provided for each kmem cache, and its name is the same as the name of the cache. For example:

```
> ::walk kmem_alloc_24
704ba010
702ba008
704ba038
702ba030
...
> ::walk thread_cache
70b38080
70aac060
705c4020
70aac1e0
```

Now you know how to iterate over the kernel memory allocator's internal data structures and examine the most important members of the kmem_cache data structure.

Detecting Memory Corruption

One of the primary debugging facilities of the allocator is that it includes algorithms to recognize data corruption quickly. When corruption is detected, the allocator immediately panics the system.

This section describes how the allocator recognizes data corruption; you must understand this to be able to debug these problems. Memory abuse typically falls into one of the following categories:

- Writing past the end of a buffer
- · Accessing uninitialized data
- · Continuing to use a freed buffer
- Corrupting kernel memory

Keep these problems in mind as you read the next three sections. They will help you to understand the allocator's design, and enable you to diagnose problems more efficiently.

Freed Buffer Checking: 0xdeadbeef

When the KMF_DEADBEEF (0x2) bit is set in the flags field of a kmem_cache, the allocator tries to make memory corruption easy to detect by writing a special pattern into all freed buffers. This pattern is

Oxdeadbeef. Since a typical region of memory contains both allocated and freed memory, sections of each kind of block will be interspersed; here is an example from the "kmem alloc 24" cache:

0x70a9add8:	deadbeef	deadbeef
0x70a9ade0:	deadbeef	deadbeef
0x70a9ade8:	deadbeef	deadbeef
0x70a9adf0:	feedface	feedface
0x70a9adf8:	70ae3260	8440c68e
0x70a9ae00:	5	4ef83
0x70a9ae08:	0	0
0x70a9ae10:	1	bbddcafe
0x70a9ae18:	feedface	139d
0x70a9ae20:	70ae3200	dlbefaed
0x70a9ae28:	deadbeef	deadbeef
0x70a9ae30:	deadbeef	deadbeef
0x70a9ae38:	deadbeef	deadbeef
0x70a9ae40:	feedface	feedface
0x70a9ae48:	70ae31a0	8440c54e

The buffer beginning at 0x70a9add8 is filled with the 0xdeadbeef pattern, which is an immediate indication that the buffer is currently free. At 0x70a9ae28 another free buffer begins; at 0x70a9ae00 an allocated buffer is located between them.

Note

You might have observed that there are some holes on this picture, and that 3 24–byte regions should occupy only 72 bytes of memory, instead of the 120 bytes shown here. This discrepancy is explained in the next section Redzone: 0xfeedface.

Redzone: 0xfeedface

The pattern <code>0xfeedface</code> appears frequently in the buffer above. This pattern is known as the "redzone" indicator. It enables the allocator (and a programmer debugging a problem) to determine if the boundaries of a buffer have been violated by "buggy" code. Following the redzone is some additional information. The contents of that data depends upon other factors (see Memory Allocation Logging). The redzone and its suffix are collectively called the *buftag* region. Figure 9–1 summarizes this information.

Figure 9.1. The Redzone

Graphic described by context.

The buftag is appended to each buffer in a cache when any of the KMF_AUDIT, KMF_DEADBEEF, or KMF_REDZONE flags are set in that buffer's cache. The contents of the buftag depend on whether KMF_AUDIT is set.

Decomposing the memory region presented above into distinct buffers is now simple:

```
0x70a9add8:
                deadbeef
                                 deadbeef
0x70a9ade0:
                deadbeef
                                 deadbeef
                                            +- User Data (free)
0x70a9ade8:
                deadbeef
                                 deadbeef
                feedface
                                 feedface -- REDZONE
0x70a9adf0:
                                 8440c68e -- Debugging Data
0x70a9adf8:
                70ae3260
0x70a9ae00:
                                 4ef83
```

```
0x70a9ae08:
                0
                                            +- User Data (allocated)
0x70a9ae10:
                1
                                bbddcafe /
0x70a9ae18:
                                139d
                                         -- REDZONE
                feedface
0x70a9ae20:
                70ae3200
                                dlbefaed -- Debugging Data
0x70a9ae28:
                deadbeef
                                deadbeef
0x70a9ae30:
                deadbeef
                                deadbeef
                                           +- User Data (free)
0x70a9ae38:
                deadbeef
                                deadbeef /
                                          -- REDZONE
                feedface
0x70a9ae40:
                                feedface
0x70a9ae48:
                70ae31a0
                                8440c54e -- Debugging Data
```

In the free buffers at 0x70a9add8 and 0x70a9ae28, the redzone is filled with 0xfeedfacefeedface. This a convenient way of determining that a buffer is free.

In the allocated buffer beginning at 0x70a9ae00, the situation is different. Recall from Allocator Basics that there are two allocation types:

- 1) The client requested memory using kmem_cache_alloc, in which case the size of the requested buffer is equal to the bufsize of the cache.
- 2) The client requested memory using kmem_alloc(9F), in which case the size of the requested buffer is less than or equal to the bufsize of the cache. For example, a request for 20 bytes will be fulfilled from the kmem_alloc_24 cache. The allocator enforces the buffer boundary by placing a marker, the *redzone byte*, immediately following the client data:

```
0x70a9ae00:
                5
                                 4ef83
0x70a9ae08:
                0
                                             +- User Data (allocated)
0x70a9ae10:
                1
                                 bbddcafe
0x70a9ae18:
                feedface
                                 139d
                                          -- REDZONE
0x70a9ae20:
                70ae3200
                                 dlbefaed -- Debugging Data
```

0xfeedface at 0x70a9ae18 is followed by a 32-bit word containing what seems to be a random value. This number is actually an encoded representation of the size of the buffer. To decode this number and find the size of the allocated buffer, use the formula:

```
size = redzone_value / 251
So, in this example,
size = 0x139d / 251 = 20 bytes.
```

This indicates that the buffer requested was of size 20 bytes. The allocator performs this decoding operation and finds that the redzone byte should be at offset 20. The redzone byte is the hex pattern 0xbb, which is present at 0x729084e4 (0x729084d0 + 0t20) as expected.

Figure 9.2. Sample kmem_alloc(9F) Buffer

This graphic depicts a sample kmem_alloc buffer. The redzone byte, uninitialized data, and debugging data are marked.

Figure 9–3 shows the general form of this memory layout.

Figure 9.3. Redzone Byte

This graphic shows the redzone byte being written after the end of the user data region. The redzone byte is determined by decoding the index.

If the allocation size is the same as the bufsize of the cache, the redzone byte overwrites the first byte of the redzone itself, as shown in Figure 9–4.

Figure 9.4. Redzone Byte at the Beginning of the Redzone

Graphic described by context.

This overwriting results in the first 32-bit word of the redzone being <code>Oxbbedface</code>, or <code>Oxfeedfabb</code> depending on the endianness of the hardware on which the system is running.

Note

Why is the allocation size encoded this way? To encode the size, the allocator uses the formula (251 * size + 1). When the size decode occurs, the integer division discards the remainder of '+1'. However, the addition of 1 is valuable because the allocator can check whether the size is valid by testing whether (size % 251 == 1). In this way, the allocator defends against corruption of the redzone byte index.

Uninitialized Data: 0xbaddcafe

You might be wondering what the suspicious <code>0xbbddcafe</code> at address <code>0x729084d4</code> was *before* the redzone byte got placed over the first byte in the word. It was <code>0xbaddcafe</code>. When the KMF_DEADBEEF flag is set in the cache, allocated but *uninitialized* memory is filled with the <code>0xbaddcafe</code> pattern. When the allocator performs an allocation, it loops across the words of the buffer and verifies that each word contains <code>0xdeadbeef</code>, then fills that word with <code>0xbaddcafe</code>.

A system can panic with a message such as:

```
panic[cpu1]/thread=e1979420: BAD TRAP: type=e (Page Fault)
rp=ef641e88 addr=baddcafe occurred in module "unix" due to an
illegal access to a user address
```

In this case, the address that caused the fault was 0xbaddcafe: the panicking thread has accessed some data that was never initialized.

Associating Panic Messages With Failures

The kernel memory allocator emits panic messages corresponding to the failure modes described earlier. For example, a system can panic with a message such as:

```
kernel memory allocator: buffer modified after being freed modification occurred at offset 0x30
```

The allocator was able to detect this case because it tried to validate that the buffer in question was filled with 0xdeadbeef. At offset 0x30, this condition was not met. Since this condition indicates memory corruption, the allocator panicked the system.

Another example failure message is:

```
kernel memory allocator: redzone violation: write past end of buffer
```

The allocator was able to detect this case because it tried to validate that the redzone byte (0xbb) was in the location it determined from the redzone size encoding. It failed to find the signature byte in the correct

location. Since this indicates memory corruption, the allocator panicked the system. Other allocator panic messages are discussed later.

Memory Allocation Logging

This section explains the logging features of the kernel memory allocator and how you can employ them to debug system crashes.

Buftag Data Integrity

As explained earlier, the second half of each buftag contains extra information about the corresponding buffer. Some of this data is debugging information, and some is data private to the allocator. While this auxiliary data can take several different forms, it is collectively known as "Buffer Control" or *bufctl* data.

However, the allocator needs to know whether a buffer's bufctl pointer is valid, since this pointer might also have been corrupted by malfunctioning code. The allocator confirms the integrity of its auxiliary pointer by storing the pointer *and* an encoded version of that pointer, and then cross-checking the two versions.

As shown in Figure 9–5, these pointers are the *bcp* (buffer control pointer) and *bxstat* (buffer control XOR status). The allocator arranges bcp and bxstat so that the expression bcp XOR bxstat equals a well-known value.

Figure 9.5. Extra Debugging Data in the Buftag

Graphic is described by context.

In the event that one or both of these pointers becomes corrupted, the allocator can easily detect such corruption and panic the system. When a buffer is *allocated*, bcp XOR bxstat = 0xalloc8ed ("allocated"). When a buffer is free, bcp XOR bxstat = 0xf4eef4ee ("freefree").

Note

You might find it helpful to re-examine the example provided in Freed Buffer Checking: 0xdeadbeef, in order to confirm that the buftag pointers shown there are consistent.

In the event that the allocator finds a corrupt buftag, it panies the system and produces a message similar to the following:

```
kernel memory allocator: boundary tag corrupted
  bcp ^ bxstat = 0xffeef4ee, should be f4eef4ee
```

Remember, if bcp is corrupt, it is still possible to retrieve its value by taking the value of bxstat XOR 0xf4eef4ee or bxstat XOR 0xall0c8ed, depending on whether the buffer is allocated or free.

The bufctl Pointer

The buffer control (bufctl) pointer contained in the buftag region can have different meanings, depending on the cache's kmem_flags. The behavior toggled by the KMF_AUDIT flag is of particular interest: when the KMF_AUDIT flag is *not* set, the kernel memory allocator allocates a kmem_bufctl_t structure for each buffer. This structure contains some minimal accounting information about each buffer. When the KMF_AUDIT flag *is* set, the allocator instead allocates a kmem_bufctl_audit_t, an extended version of the kmem_bufctl_t.

This section presumes the KMF_AUDIT flag is set. For caches that do not have this bit set, the amount of available debugging information is reduced.

The kmem_bufctl_audit_t (bufctl_audit for short) contains additional information about the last transaction that occurred on this buffer. The following example shows how to apply the bufctl_audit macro to examine an audit record. The buffer shown is the example buffer used in Detecting Memory Corruption:

Using the techniques presented above, it is easy to see that 0x70ae3200 points to the bufctl_audit record: it is the first pointer following the redzone. To examine the bufctl_audit record it points to, apply the bufctl_audit macro:

> 0x70ae3200\$ <bufctl_audit< th=""></bufctl_audit<>			
0x70ae3200:	next	addr	slab
	70378000	70a9ae00	707c86a0
0x70ae320c:	cache	timestamp	thread
	70039928	e1bd0e26afe	70aac4e0
0x70ae321c:	lastlog	contents	stackdepth
	7011c7c0	7018a0b0	4
0x70ae3228:			
	kmem_zalloc+0x3	0	
	pid_assign+8		
	getproc+0x68		
	cfork+0x60		

The 'addr' field is the address of the buffer corresponding to this bufctl_audit record. This is the original address: 0x70a9ae00. The 'cache' field points at the kmem_cache that allocated this buffer. You can use the ::kmem cache dcmd to examine it as follows:

```
> 0x70039928::kmem_cache

ADDR NAME FLAG CFLAG BUFSIZE BUFTOTL

70039928 kmem alloc 24 020f 000000 24 612
```

The 'timestamp' field represents the time this transaction occurred. This time is expressed in the same manner as gethrtime(3C).

'thread' is a pointer to the thread that performed the last transaction on this buffer. The 'lastlog' and 'contents' pointers point to locations in the allocator's *transaction logs*. These logs are discussed in detail in Allocator Logging Facility.

Typically, the most useful piece of information provided by bufctl_audit is the stack trace recorded at the point at which the transaction took place. In this case, the transaction was an allocation called as part of executing fork(2).

Advanced Memory Analysis

This section describes facilities for performing advanced memory analysis, including locating memory leaks and sources of data corruption.

Finding Memory Leaks

The ::findleaks dcmd provides powerful and efficient detection of memory leaks in kernel crash dumps where the full set of kmem debug features has been enabled. The first execution of ::findleaks processes the dump for memory leaks (this can take a few minutes), and then coalesces the leaks by the allocation stack trace. The findleaks report shows a bufctl address and the topmost stack frame for each memory leak that was identified:

Using the bufctl pointers, you can obtain the complete stack backtrace of the allocation by applying the bufctl_audit macro:

```
> 70d3b1a0$<bufctl_audit</pre>
0x70d3b1a0:
                next
                                  addr
                                                   slab
                70a049c0
                                  70d03b28
                                                   70bb7480
0x70d3b1ac:
                cache
                                  timestamp
                                                   thread
                7003a028
                                  13f7cf63b3
                                                   70b38380
0x70d3b1bc:
                lastlog
                                                   stackdepth
                                  contents
                 700d6e60
0x70d3b1c8:
                kmem_alloc+0x30
                 sigaddq+0x108
                sigsendproc+0x210
                 sigqkill+0x90
                kill+0x28
```

The programmer can usually use the bufctl_audit information and the allocation stack trace to quickly track down the code path that leaks the given buffer.

Finding References to Data

When trying to diagnose a memory corruption problem, you should know what other kernel entities hold a copy of a particular pointer. This is important because it can reveal which thread accessed a data structure after it was freed. It can also make it easier to understand what kernel entities are sharing knowledge of a particular (valid) data item. The ::whatis and ::kgrep dcmds can be used to answer these questions. You can apply ::whatis to a value of interest:

```
> 0x705d8640::whatis
705d8640 is 705d8640+0, allocated from streams_mblk
```

In this case, 0x705d8640 is revealed to be a pointer to a STREAMS mblk structure. To see the entire allocation tree, use ::whatis a instead:

```
> 0x705d8640::whatis -a
705d8640 is 705d8640+0, allocated from streams_mblk
705d8640 is 705d8000+640, allocated from kmem_va_8192
```

```
705d8640 is 705d8000+640 from kmem_default vmem arena 705d8640 is 705d2000+2640 from kmem_va vmem arena 705d8640 is 705d2000+2640 from heap vmem arena
```

This reveals that the allocation also appears in the kmem_va_8192 cache--a kmem cache that is fronting the kmem va vmem arena. It also shows the full stack of vmem allocations.

The complete list of kmem caches and vmem arenas is displayed by the ::kmastat dcmd. You can use ::kgrep to locate other kernel addresses that contain a pointer to this mblk. This illustrates the hierarchical nature of memory allocations in the system; in general, you can determine the type of object referred to by the given address from the name of the most specific kmem cache.

```
> 0x705d8640::kgrep
400a3720
70580d24
7069d7f0
706a37ec
706add34
```

and investigate them by applying ::whatis again:

```
> 400a3720::whatis
400a3720 is in thread 7095b240's stack
> 706add34::whatis
706add34 is 706add20+14, allocated from streams_dblk_120
```

Here one pointer is located on the stack of a known kernel thread, and another is the mblk pointer inside of the corresponding STREAMS dblk structure.

Finding Corrupt Buffers With ::kmem_verify

MDB's ::kmem_verify dcmd implements most of the same checks that the kmem allocator does at runtime. ::kmem_verify can be invoked in order to scan every kmem cache with appropriate kmem_flags, or to examine a particular cache.

Here is an example of using ::kmem verify to isolate a problem:

```
> ::kmem verify
Cache Name
                                 Addr
                                          Cache Integrity
kmem alloc 8
                                 70039428 clean
kmem_alloc_16
                                 700396a8 clean
kmem alloc 24
                                 70039928 1 corrupt buffer
kmem alloc 32
                                 70039ba8 clean
kmem alloc 40
                                 7003a028 clean
kmem_alloc_48
                                 7003a2a8 clean
```

It is easy to see here that the kmem_alloc_24 cache contains what ::kmem_verify believes to be a problem. With an explicit cache argument, the ::kmem_verify dcmd provides more detailed information about the problem:

```
> 70039928::kmem_verify
Summary for cache 'kmem_alloc_24'
buffer 702babc0 (free) seems corrupted, at 702babc0
```

The next step is to examine the buffer which ::kmem_verify believes to be corrupt:

The reason that ::kmem_verify flagged this buffer is now clear: The first word in the buffer (at 0x702babc0) should probably be filled with the 0xdeadbeef pattern, not with a 0. At this point, examining the bufctl_audit for this buffer might yield clues about what code recently wrote to the buffer, indicating where and when it was freed.

Another useful technique in this situation is to use ::kgrep to search the address space for references to address 0x702babc0, in order to discover what threads or data structures are still holding references to this freed data.

Allocator Logging Facility

When KMF_AUDIT is set for a cache, the kernel memory allocator maintains a log that records the recent history of its activity. This *transaction log* records bufctl_audit records. If the KMF_AUDIT and the KMF_CONTENTS flags are both set, the allocator generates a *contents log* that records portions of the actual contents of allocated and freed buffers. The structure and use of the contents log is outside the scope of this document. The transaction log is discussed in this section.

MDB provides several facilities for displaying the transaction log. The simplest is ::walk kmem_log, which prints out the transaction in the log as a series of bufctl_audit_t pointers:

```
> ::walk kmem_log
70128340
701282e0
70128280
70128220
701281c0
> 70128340$<bufctl audit
0 \times 70128340:
                 next
                                  addr
                                                    slab
                 70ac1d40
                                  70bc4ea8
                                                   70bb7c00
0x7012834c:
                 cache
                                  timestamp
                                                    thread
                 70039428
                                                   70aacde0
                                  e1bd7abe721
0x7012835c:
                 lastlog
                                                   stackdepth
                                  contents
                 701282e0
                                  7018f340
0x70128368:
                 kmem_cache_free+0x24
                 nfs3_sync+0x3c
                 vfs_sync+0x84
                 syssync+4
```

A more elegant way to view the entire transaction log is by using the ::kmem_log command:

```
> ::kmem_log
CPU ADDR BUFADDR TIMESTAMP THREAD
    0 70128340 70bc4ea8 elbd7abe721 70aacde0
    0 701282e0 70bc4ea8 elbd7aa86fa 70aacde0
```

```
0 70128280 70bc4ea8
                         e1bd7aa27dd 70aacde0
0 70128220 70bc4ea8
                         elbd7a98a6e 70aacde0
0 701281c0 70d03738
                         e1bd7a8e3e0 70aacde0
                         e1bd78035ad 70aacde0
0 70127140 70cf78a0
0 701270e0 709cf6c0
                         e1bd6d2573a 40033e60
0 70127080 70cedf20
                         e1bd6d1e984 40033e60
0 70127020 70b09578
                         e1bd5fc1791 40033e60
0 70126fc0 70cf78a0
                         elbd5fb6b5a 40033e60
0 70126f60 705ed388
                         e1bd5fb080d 40033e60
0 70126f00 705ed388
                         e1bd551ff73 70aacde0
```

The output of ::kmem_log is sorted in descending order by timestamp. The ADDR column is the bufctl_audit structure corresponding to that transaction; BUFADDR points to the actual buffer.

These figures represent *transactions* on buffers (both allocations and frees). When a particular buffer is corrupted, it can be helpful to locate that buffer in the transaction log, then determine in which other transactions the transacting thread was involved. This can help to assemble a picture of the sequence of events that occurred prior to and after the allocation (or free) of a buffer.

You can employ the ::bufctl command to filter the output of walking the transaction log. The ::bufctl -a command filters the buffers in the transaction log by buffer address. This example filters on buffer 0x70b09578:

```
> ::walk kmem log | ::bufctl -a 0x70b09578
ADDR
        BUFADDR
                  TIMESTAMP
                             THREAD
                                       CALLER
70127020 70b09578 e1bd5fc1791 40033e60 biodone+0x108
70126e40 70b09578 e1bd55062da 70aacde0 pageio setup+0x268
70126de0 70b09578 e1bd52b2317 40033e60 biodone+0x108
70126c00 70b09578 e1bd497ee8e 70aacde0 pageio_setup+0x268
70120480 70b09578 e1bd21c5e2a 70aacde0 elfexec+0x9f0
70120060 70b09578 e1bd20f5ab5 70aacde0 getelfhead+0x100
7011ef20 70b09578 elbdle9aldd 70aacde0 ufs_getpage_miss+0x354
7011d720 70b09578 e1bd1170dc4 70aacde0 pageio setup+0x268
70117d80 70b09578 e1bcff6ff27 70bc2480 elfexec+0x9f0
70117960 70b09578 elbcfea4a9f 70bc2480 getelfhead+0x100
```

This example illustrates that a particular buffer can be used in numerous transactions.

Note

Remember that the kmem transaction log is an incomplete record of the transactions made by the kernel memory allocator. Older entries in the log are evicted as needed in order to keep the size of the log constant.

The ::allocdby and ::freedby dcmds provide a convenient way to summarize transactions associated with a particular thread. Here is an example of listing the recent allocations performed by thread 0x70aacde0:

```
> 0x70aacde0::allocdby
BUFCTL TIMESTAMP CALLER
70d4d8c0 eledb14511a allocb+0x88
70d4e8a0 eledb142472 dblk_constructor+0xc
70d4a240 eledb13dd4f allocb+0x88
```

70d4e840	eledb13aeec	dblk_constructor+0xc
70d4d860	e1ed8344071	allocb+0x88
70d4e7e0	e1ed8342536	dblk_constructor+0xc
70d4a1e0	e1ed82b3a3c	allocb+0x88
70a53f80	e1ed82b0b91	dblk_constructor+0xc
70d4d800	e1e9b663b92	allocb+0x88

By examining bufctl_audit records, you can understand the recent activities of a particular thread.

Chapter 10. Module Programming API

This chapter describes the structures and functions contained in the MDB debugger module API. The header file <sys/mdb_modapi.h> contains prototypes for these functions, and the SUNWmdbdm package provides source code for an example module in the directory /usr/demo/mdb.

Debugger Module Linkage

mdb init

```
const mdb_modinfo_t *_mdb_init(void);
```

Each debugger module is required to provide, for linkage and identification purposes, a function named _mdb_init. This function returns a pointer to a persistent (that is, not declared as an automatic variable) mdb_modinfo_t structure, as defined in <sys/mdb_modapi.h>:

The *mi_dvers* member is used to identify the API version number, and should always be set to MDB_API_VERSION. The current version number is therefore compiled into each debugger module, allowing the debugger to identify and verify the application binary interface used by the module. The debugger does not load modules that are compiled for an API version that is more recent than the debugger itself.

The mi_dcmds and mi_walkers members, if not NULL, point to arrays of dcmd and walker definition structures, respectively. Each array must be terminated by a NULL element. These dcmds and walkers are installed and registered with the debugger as part of the module loading process. The debugger will refuse to load the module if one or more dcmds or walkers are defined improperly or if they have conflicting or invalid names. Dcmd and walker names are prohibited from containing characters that have special meaning to the debugger, such as quotation marks and parentheses.

The module can also execute code in _mdb_init using the module API to determine if it is appropriate to load. For example, a module can only be appropriate for a particular target if certain symbols are present. If these symbols are not found, the module can return NULL from the _mdb_init function. In this case, the debugger will refuse to load the module and an appropriate error message is printed.

_mdb_fini

```
void _mdb_fini(void);
```

If the module performs certain tasks prior to unloading, such as freeing persistent memory previously allocated with mdb_alloc, it can declare a function named _mdb_fini for this purpose. This function is not required by the debugger. If declared, it is called once prior to unloading the module. Modules are unloaded when the user requests that the debugger terminate or when the user explicitly unloads a module using the ::unload built-in dcmd.

Dcmd Definitions

```
int dcmd(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv);
```

A dcmd is implemented with a function similar to the dcmd declaration. This function receives four arguments and returns an integer status. The function arguments are:

addrCurrent address, also called dot. At the start of the dcmd, this address corresponds to the value of the dot "." variable in the debugger.

flagInteger containing the logical OR of one or more of the following flags:

DCMDAADARBAEAddress was specified to the left of ::dcmd.

DCMD**Thede**md was invoked in a loop using the ,count syntax, or the dcmd was invoked in a loop by a pipeline.

DCMD**This One Notice that the design of the demd function corresponds to the first loop or pipeline invocation.**

DCMDThe domd was invoked with input from a pipeline.

DCMDThe domowas invoked with output set to a pipeline.

As a convenience, the DCMD_HDRSPEC macro is provided to allow a dcmd to test its flags to determine if it should print a header line (that is, it was not invoked as part of a loop, or it was invoked as the first iteration of a loop or pipeline).

argcNumber of arguments in the argv array.

argvArray of arguments specified to the right of ::dcmd on the command line. These arguments can be either strings or integer values.

The dcmd function is expected to return one of the following integer values, defined in <sys/mdb_modapi.h>.

DCMD<u>Tb</u>rdcmd completed successfully.

DCMDTER cmd failed for some reason.

DCMD**The Acte**d failed because invalid arguments were specified. When this value is returned, the dcmd usage message (described below) prints automatically.

DCMDTNE NEXt dcmd definition (if one is present) is automatically invoked with the same arguments.

DCMD_TABORTIC failed, and the current loop or pipeline should be aborted. This is like DCMD_ERR, but indicates that no further progress is possible in the current loop or pipe.

Each dcmd consists of a function defined according to the example dcmd prototype, and a corresponding mdb_dcmd_t structure, as defined in <sys/mdb_modapi.h>. This structure consists of the following fields:

const char *dc_name The string name of the dcmd, without the leading "::". The name cannot contain any of the MDB meta-characters, such as \$ or `.

const char *dc_usage An optional usage string for the dcmd, to be printed when the dcmd returns DCMD_USAGE. For example, if the dcmd accepts options a and b, dc_usage might be specified as "[ab]". If the dcmd accepts no arguments, dc_usage can be set to NULL. If the usage string begins

with ":", this is shorthand for indicating that the dcmd requires an explicit address (that is, it requires DCMD_ADDRSPEC to be set in its

flags parameter). If the usage string begins with "?", this indicates that the dcmd optionally accepts an address. These hints modify the usage message accordingly.

const char *dc_descr A mandatory description string, briefly explaining the purpose of the dcmd. This string should consist of only a single line of text.

mdb_dcmd_f *dc_funcp A pointer to the function that will be called to execute the dcmd.

void (*dc_help)(void) An optional function pointer to a help function for the dcmd. If this pointer is not NULL, this function will be called when the user executes ::help dcmd. This function can use mdb_printf to display further

information or examples.

Walker Definitions

```
int walk_init(mdb_walk_state_t *wsp);
int walk_step(mdb_walk_state_t *wsp);
void walk_fini(mdb_walk_state_t *wsp);
```

A walker is composed of three functions, init, step, and fini, which are defined according to the example prototypes above. A walker is invoked by the debugger when one of the walk functions (such as mdb_walk) is called, or when the user executes the ::walk built-in dcmd. When the walk begins, MDB calls the walker's init function, passing it the address of a new mdb_walk_state_t structure, as defined in <sys/mdb modapi.h>:

A separate mdb_walk_state_t is created for each walk, so that multiple instances of the same walker can be active simultaneously. The state structure contains the callback the walker should invoke at each step (walk_callback), and the private data for the callback (walk_cbdata), as specified to mdb_walk, for example. The walk_cbdata pointer is opaque to the walker: it must not modify or dereference this value, nor can it assume it is a pointer to valid memory.

The starting address for the walk is stored in walk_addr. This is either NULL if mdb_walk was called, or the address parameter specified to mdb_pwalk. If the ::walk built-in was used, walk_addr will be non-NULL if an explicit address was specified on the left-hand side of ::walk. A walk with a starting address of NULL is referred to as global. A walk with an explicit non-NULL starting address is referred to as local.

The walk_data and walk_arg fields are provided for use as private storage for the walker. Complex walkers might need to allocate an auxiliary state structure and set walk_data to point to this structure. Each time a walk is initiated, walk_arg is initialized to the value of the walk_init_arg member of the corresponding walker's mdb_walker_t structure.

In some cases, it is useful to have several walkers share the same init, step, and fini routines. For example, the MDB genunix module provides walkers for each kernel memory cache. These share the same init,

step, and fini functions, and use the walk_init_arg member of the mdb_walker_t to specify the address of the appropriate cache as the walk arg.

If the walker calls mdb_layered_walk to instantiate an underlying layer, then the underlying layer will reset walk_addr and walk_layer prior to each call to the walker's step function. The underlying layer sets walk_addr to the target virtual address of the underlying object, and set walk_layer to point to the walker's local copy of the underlying object. For more information on layered walks, refer to the discussion of mdb_layered_walk below.

The walker init and step functions are expected to return one of the following status values:

- WALK_NEXT, MDB invokes the walk step function. When the walk step function returns WALK_NEXT, this indicates that MDB should call the step function again.
- WALK_DONE can be returned by either the step function to indicate that the walk is complete, or by the init function to indicate that no steps are needed (for example, if the given data structure is empty).
- WALK_TERN walk has terminated due to an error. If WALK_ERR is returned by the init function, mdb_walk (or any of its counterparts) returns -1 to indicate that the walker failed to initialize. If WALK_ERR is returned by the step function, the walk terminates but mdb_walk returns success.

The walk_callback is also expected to return one of the values above. Therefore, the walk step function's job is to determine the address of the next object, read in a local copy of this object, call the walk_callback function, then return its status. The step function can also return WALK_DONE or WALK_ERR without invoking the callback if the walk is complete or if an error occurred.

The walker itself is defined using the mdb_walker_t structure, defined in:

The walk_name and walk_descr fields should be initialized to point to strings containing the name and a brief description of the walker, respectively. A walker is required to have a non-NULL name and description, and the name cannot contain any of the MDB meta-characters. The description string is printed by the ::walkers and ::dmods built-in dcmds.

The walk_init, walk_step, and walk_fini members refer to the walk functions themselves, as described earlier. The walk_init and walk_fini members can be set to NULL to indicate that no special initialization or cleanup actions need to be taken. The walk_step member cannot be set to NULL. The walk_init_arg member is used to initialize the walk_arg member of each new mdb_walk_state_t created for the given walker, as described earlier. Figure 10–1 shows a flowchart for the algorithm of a typical walker.

Figure 10.1. Sample Walker

Graphic is described by context.

The walker is designed to iterate over the list of proc_t structures in the kernel. The head of the list is stored in the global practive variable, and each element's p_next pointer points to the next proc_t in

the list. The list is terminated with a NULL pointer. In the walker's init routine, the practive symbol is located using mdb_lookup_by_name step (1), and its value is copied into the mdb_walk_state_t pointed to by wsp.

In the walker's step function, the next proc_t structure in the list is copied into the debugger's address space using mdb_vread step (2), the callback function is invoked with a pointer to this local copy, step (3), and then the mdb_walk_state_t is updated with the address of the proc_t structure for the next iteration. This update corresponds to following the pointer, step (4), to the next element in the list.

These steps demonstrate the structure of a typical walker: the init routine locates the global information for a particular data structure, the step function reads in a local copy of the next data item and passes it to the callback function, and the address of the next element is read. Finally, when the walk terminates, the fini function frees any private storage.

API Functions

mdb_pwalk

Initiate a local walk starting at <code>addr</code> using the walker specified by <code>name</code>, and invoke the callback function <code>func</code> at each step. If <code>addr</code> is NULL, a global walk is performed (that is, the mdb_pwalk invocation is equivalent to the identical call to mdb_walk without the trailing <code>addr</code> parameter). This function returns 0 for success, or -1 for error. The mdb_pwalk function fails if the walker itself returns a fatal error, or if the specified walker name is not known to the debugger. The walker name may be scoped using the backquote (`) operator if there are naming conflicts. The <code>data</code> parameter is an opaque argument that has meaning only to the caller; it is passed back to <code>func</code> at each step of the walk.

mdb_walk

```
int mdb_walk(const char *name, mdb_walk_cb_t func, void *data);
```

Initiate a global walk starting at addr using the walker specified by name, and invoke the callback function func at each step. This function returns 0 for success, or -1 for error. The mdb_walk function fails if the walker itself returns a fatal error, or if the specified walker name is not known to the debugger. The walker name can be scoped using the backquote (`) operator if there are naming conflicts. The data parameter is an opaque argument that has meaning only to the caller; it is passed back to func at each step of the walk.

mdb_pwalk_dcmd

```
int mdb_pwalk_dcmd(const char *wname, const char *dcname, int argc,
    const mdb_arg_t *argv, uintptr_t addr);
```

Initiate a local walk starting at addr using the walker specified by wname, and invoke the dcmd specified by dcname with the specified argc and argv at each step. This function returns 0 for success, or -1 for error. The function fails if the walker itself returns a fatal error, if the specified walker name or dcmd name is not known to the debugger, or if the dcmd itself returns DCMD_ABORT or DCMD_USAGE to the walker. The walker name and dcmd name can each be scoped using the backquote (`) operator if there are naming conflicts. When invoked from mdb_pwalk_dcmd, the dcmd will have the DCMD_LOOP and DCMD_ADDRSPEC bits set in its flags parameter, and the first call will have DCMD_LOOPFIRST set.

mdb walk dcmd

```
int mdb_walk_dcmd(const char *wname, const char *dcname, int argc,
    const mdb_arg_t *argv);
```

Initiate a global walk using the walker specified by wname, and invoke the dcmd specified by dcname with the specified argc and argv at each step. This function returns 0 for success, or -1 for error. The function fails if the walker itself returns a fatal error, if the specified walker name or dcmd name is not known to the debugger, or if the dcmd itself returns DCMD_ABORT or DCMD_USAGE to the walker. The walker name and dcmd name can each be scoped using the backquote (`) operator if there are naming conflicts. When invoked from mdb_walk_dcmd, the dcmd will have the DCMD_LOOP and DCMD_ADDRSPEC bits set in its flags parameter, and the first call will have DCMD_LOOPFIRST set.

mdb call dcmd

```
int mdb_call_dcmd(const char *name, uintptr_t addr, uint_t flags,
  int argc, const mdb_arg_t *argv);
```

Invoke the specified dcmd name with the given parameters. The dot variable is reset to addr, and addr, flags, argc, and argv are passed to the dcmd. The function returns 0 for success, or -1 for error. The function fails if the dcmd returns DCMD_ERR, DCMD_ABORT, or DCMD_USAGE, or if the specified dcmd name is not known to the debugger. The dcmd name can be scoped using the backquote (`) operator if there are naming conflicts.

mdb_layered_walk

```
int mdb_layered_walk(const char *name, mdb_walk_state_t *wsp);
```

Layer the walk denoted by wsp on top of a walk initiated using the specified walker name. The name can be scoped using the backquote (`) operator if there are naming conflicts. Layered walks can be used, for example, to facilitate constructing walkers for data structures that are embedded in other data structures.

For example, suppose that each CPU structure in the kernel contains a pointer to an embedded structure. To write a walker for the embedded structure type, you could replicate the code to iterate over CPU structures and dereference the appropriate member of each CPU structure, or you could layer the embedded structure's walker on top of the existing CPU walker.

The mdb_layered_walk function is used from within a walker's init routine to add a new layer to the current walk. The underlying layer is initialized as part of the call to mdb_layered_walk. The calling walk routine passes in a pointer to its current walk state; this state is used to construct the layered walk. Each layered walk is cleaned up after the caller's walk fini function is called. If more than one layer is added to a walk, the caller's walk step function will step through each element returned by the first layer, then the second layer, and so forth.

The mdb_layered_walk function returns 0 for success, or -1 for error. The function fails if the specified walker name is not known to the debugger, if the wsp pointer is not a valid, active walk state pointer, if the layered walker itself fails to initialize, or if the caller attempts to layer the walker on top of itself.

mdb_add_walker

```
int mdb_add_walker(const mdb_walker_t *w);
```

Register a new walker with the debugger. The walker is added to the module's namespace, and to the debugger's global namespace according to the name resolution rules described in dcmd and Walker

Name Resolution. This function returns 0 for success, or -1 for error if the given walker name is already registered by this module, or if the walker structure w is improperly constructed. The information in the mdb_walker_t w is copied to internal debugger structures, so the caller can reuse or free this structure after the call to mdb_add_walker.

mdb_remove_walker

```
int mdb_remove_walker(const char *name);
```

Remove the walker with the specified name. This function returns 0 for success, or -1 for error. The walker is removed from the current module's namespace. The function fails if the walker name is unknown, or is registered only in another module's namespace. The mdb_remove_walker function can be used to remove walkers that were added dynamically using mdb_add_walker, or walkers that were added statically as part of the module's linkage structure. The scoping operator cannot be used in the walker name; it is not legal for the caller of mdb_remove_walker to attempt to remove a walker exported by a different module.

mdb vread and mdb vwrite

```
ssize_t mdb_vread(void *buf, size_t nbytes, uintptr_t addr);
ssize_t mdb_vwrite(const void *buf, size_t nbytes, uintptr_t addr);
```

These functions provide the ability to read and write data from a given target virtual address, specified by the addr parameter. The mdb_vread function returns nbytes for success, or -1 for error; if a read is truncated because only a portion of the data can be read from the specified address, -1 is returned. The mdb_vwrite function returns the number of bytes actually written upon success; -1 is returned upon error.

mdb_fread and mdb_fwrite

```
ssize_t mdb_fread(void *buf, size_t nbytes, uintptr_t addr);
ssize_t mdb_fwrite(const void *buf, size_t nbytes, uintptr_t addr);
```

These functions provide the ability to read and write data from the object file location corresponding to the given target virtual address, specified by the addr parameter. The mdb_fread function returns nbytes for success, or -1 for error; if a read is truncated because only a portion of the data can be read from the specified address, -1 is returned. The mdb_fwrite function returns the number of bytes actually written upon success; -1 is returned upon error.

mdb_pread and mdb_pwrite

```
ssize_t mdb_pread(void *buf, size_t nbytes, uint64_t addr);
ssize_t mdb_pwrite(const void *buf, size_t nbytes, uint64_t addr);
```

These functions provide the ability to read and write data from a given target physical address, specified by the addr parameter. The mdb_pread function returns nbytes for success, or -1 for error; if a read is truncated because only a portion of the data can be read from the specified address, -1 is returned. The mdb_pwrite function returns the number of bytes actually written upon success; -1 is returned upon error.

mdb_readstr

```
ssize_t mdb_readstr(char *s, size_t nbytes, uintptr_t addr);
```

The mdb_readstr function reads a null-terminated C string beginning at the target virtual address addr into the buffer addressed by s. The size of the buffer is specified by nbytes. If the string is longer than can fit in the buffer, the string is truncated to the buffer size and a null byte is stored at s[nbytes - 1]. The length of the string stored in s (not including the terminating null byte) is returned upon success; otherwise -1 is returned to indicate an error.

mdb writestr

```
ssize_t mdb_writestr(const char *s, uintptr_t addr);
```

The mdb_writestr function writes a null-terminated C string from s (including the trailing null byte) to the target's virtual address space at the address specified by addr. The number of bytes written (not including the terminating null byte) is returned upon success; otherwise, -1 is returned to indicate an error.

mdb_readsym

```
ssize_t mdb_readsym(void *buf, size_t nbytes, const char *name);
```

mdb_readsym is similar to mdb_vread, except that the virtual address at which reading begins is obtained from the value of the symbol specified by *name*. If no symbol by that name is found or a read error occurs, -1 is returned; otherwise *nbytes* is returned for success.

The caller can first look up the symbol separately if it is necessary to distinguish between symbol lookup failure and read failure. The primary executable's symbol table is used for the symbol lookup; if the symbol resides in another symbol table, you must first apply mdb_lookup_by_obj, then mdb_vread.

mdb_writesym

```
ssize_t mdb_writesym(const void *buf, size_t nbytes, const char *name);
```

mdb_writesym is identical to mdb_vwrite, except that the virtual address at which writing begins is obtained from the value of the symbol specified by name. If no symbol by that name is found, -1 is returned. Otherwise, the number of bytes successfully written is returned on success, and -1 is returned on error. The primary executable's symbol table is used for the symbol lookup; if the symbol resides in another symbol table, you must first apply mdb_lookup_by_obj, then mdb_vwrite.

mdb_readvar and mdb_writevar

```
ssize_t mdb_readvar(void *buf, const char *name);
ssize_t mdb_writevar(const void *buf, const char *name);
```

mdb_readvar is similar to mdb_vread, except that the virtual address at which reading begins and the number of bytes to read are obtained from the value and size of the symbol specified by name. If no symbol by that name is found, -1 is returned. The symbol size (the number of bytes read) is returned on success; -1 is returned on error. This is useful for reading well-known variables whose sizes are fixed. For example:

```
int hz; /* system clock rate */
mdb_readvar(&hz, "hz");
```

The caller can first look up the symbol separately if it is necessary to distinguish between symbol lookup failure and read failure. The caller must also carefully check the definition of the symbol of interest in order to make sure that the local declaration is the exact same type as the target's definition. For example, if the caller declares an int, and the symbol of interest is actually a long, and the debugger is examining a 64-bit kernel target, mdb_readvar copies back 8 bytes to the caller's buffer, corrupting the 4 bytes following the storage for the int.

mdb_writevar is identical to mdb_vwrite, except that the virtual address at which writing begins and the number of bytes to write are obtained from the value and size of the symbol specified by name. If no symbol by that name is found, -1 is returned. Otherwise, the number of bytes successfully written is returned on success, and -1 is returned on error.

For both functions, the primary executable's symbol table is used for the symbol lookup; if the symbol resides in another symbol table, you must first apply mdb_lookup_by_obj, then mdb_vread or mdb_vwrite.

mdb_lookup_by_name and mdb_lookup_by_obj

```
int mdb_lookup_by_name(const char *name, GElf_Sym *sym);
int mdb lookup by obj(const char *object, const char *name, GElf Sym *sym);
```

Look up the specified symbol name and copy the ELF symbol information into the *GE1f_Sym* pointed to by *sym*. If the symbol is found, the function returns 0; otherwise, -1 is returned. The *name* parameter specifies the symbol name. The *object* parameter tells the debugger where to look for the symbol. For the mdb_lookup_by_name function, the object file defaults to MDB_OBJ_EXEC. For mdb_lookup_by_obj, the object name should be one of the following:

- MDB_DBok_inxterexecutable's symbol table (.symtab section). For kernel crash dumps, this corresponds to the symbol table from the unix.X file or from /dev/ksyms.
- MDB_DBOK_RTIME runtime link-editor's symbol table. For kernel crash dumps, this corresponds to the symbol table for the krtld module.
- MDB_DBGk_invillerown symbol tables. For kernel crash dumps, this includes the .symtab and .dynsym sections from the unix.X file or /dev/ksyms, as well as per-module symbol tables if these have been processed.
- objetithe name of a particular load object is explicitly specified, the search is restricted to the symbol table of this object. The object can be named according to the naming convention for load objects described in Symbol Name Resolution.

mdb lookup by addr

Locate the symbol corresponding to the specified address and copy the ELF symbol information into the *GElf_Sym* pointed to by *sym* and the symbol name into the character array addressed by *buf*. If a corresponding symbol is found, the function returns 0; otherwise -1 is returned.

The flag parameter specifies the lookup mode and should be one of the following:

- MDB_AYMWFfuzzy matching to take place, based on the current symbol distance setting. The symbol distance can be controlled using the ::set s built-in. If an explicit symbol distance has been set (absolute mode), the address can match a symbol if the distance from the symbol's value to the address does not exceed the absolute symbol distance. If smart mode is enabled (symbol distance = 0), then the address can match the symbol if it is in the range [symbol value, symbol value + symbol size).
- MDB_**Distall EXACTE** y matching. The symbol can match only the address if the symbol value exactly equals the specified address.

If a symbol match occurs, the name of the symbol is copied into the *buf* supplied by the caller. The *len* parameter specifies the length of this buffer in bytes. The caller's *buf* should be at least of size MDB_SYM_NAMLEN bytes. The debugger copies the name to this buffer and appends a trailing null byte. If the name length exceeds the length of the buffer, the name is truncated but always includes a trailing null byte.

mdb_getopts

```
int mdb_getopts(int argc, const mdb_arg_t *argv, ...);
```

Parse and process options and option arguments from the specified argument array (argv). The argc parameter denotes the length of the argument array. This function processes each argument in order, and stops and returns the array index of the first argument that could not be processed. If all arguments are processed successfully, argc is returned.

Following the argc and argv parameters, the mdb_getopts function accepts a variable list of arguments describing the options that are expected to appear in the argv array. Each option is described by an option letter (char argument), an option type (uint_t argument), and one or two additional arguments, as shown in the table below. The list of option arguments is terminated with a NULL argument. The type should be one of one of the following:

```
char c, uint_t type, uint_t bits, uint_t *p
```

If type is MDB_OPT_SETBITS and option c is detected in the argv list, the debugger will OR bits into the integer referenced by pointer p.

MDB_**DET** option is described by these parameters:

```
char c, uint_t type, uint_t bits, uint_t *p
```

If type is MDB_OPT_CLRBITS and option c is detected in the argv list, the debugger clears bits from the integer referenced by pointer p.

MDB_**TRE** accepts a string argument. The option is described by these parameters:

```
char c, uint_t type, const char **p
```

If type is MDB_OPT_STR and option c is detected in the argv list, the debugger stores a pointer to the string argument following c in the pointer referenced by p.

MDB_**OPETODISMIZECTE**pts a uintptr_t argument. The option is described by these parameters:

```
char c, uint_t type, uintptr_t *p
```

If type is MDB_OPT_UINTPTR and option c is detected in the argv list, the debugger stores the integer argument following c in the uintptr_t referenced by p.

MDB_**OBETOINTER TEDISET** intptr_t argument. The option is described by these parameters:

```
char c, uint_t type, boolean_t *flag, uintptr_t *p
```

If type is MDB_OPT_UINTPTR_SET and option c is detected in the argv list, the debugger stores the value '1' (TRUE) into the boolean_t referenced by flag, and the integer argument following c in the uintptr_t referenced by p.

MDB_**DREONIZE**4epts a uint 64_t argument. The option is described by these parameters:

```
char c, uint_t type, uint64_t *p
```

If type is MDB_OPT_UINT64 and option c is detected in the argv list, the debugger stores the integer argument following c in the uint64_t referenced by p.

For example, the following source code:

demonstrates how mdb_getopts might be used in a dcmd to accept a boolean option "v" that sets the opt_v variable to TRUE, and an option "s" that accepts a string argument that is stored in the opt_s variable. The mdb_getopts function also automatically issues warning messages if it detects an invalid option letter or missing option argument before returning to the caller. The storage for argument strings and the argv array is automatically garbage-collected by the debugger upon completion of the dcmd.

mdb_strtoull

```
u_longlong_t mdb_strtoull(const char *s);
```

Convert the specified string s to an unsigned long long representation. This function is intended for use in processing and converting string arguments in situations where mdb_getopts is not appropriate. If the string argument cannot be converted to a valid integer representation, the function fails by printing an appropriate error message and aborting the dcmd. Therefore, error checking code is not required. The string can be prefixed with any of the valid base specifiers (0i, 0I, 0o, 0O, 0t, 0T, 0x, or 0X); otherwise, it is interpreted using the default base. The function will fail and abort the dcmd if any of the characters in s are not appropriate for the base, or if integer overflow occurs.

mdb_alloc, mdb_zalloc and mdb_free

```
void *mdb_alloc(size_t size, uint_t flags);
void *mdb_zalloc(size_t size, uint_t flags);
void mdb_free(void *buf, size_t size);
```

mdb_alloc allocates size bytes of debugger memory and returns a pointer to the allocated memory. The allocated memory is at least double-word aligned, so it can hold any C data structure. No greater alignment can be assumed. The flags parameter should be the bitwise OR of one or more of the following values:

UM_NIGNIFIED ent memory to fulfill the request is not immediately available, return NULL to indicate failure. The caller must check for NULL and handle this case appropriately.

- UM_SIGNATICE
 UM_SIGNATICE
 Item ficient memory to fulfill the request is not immediately available, sleep until such time as the request can be fulfilled. As a result, UM_SLEEP allocations are guaranteed to succeed. The caller need not check for a NULL return value.
- UM_GGarbage-collect allocation automatically at the end of this debugger command. The caller should not subsequently call mdb_free on this block, as the debugger will take care of de-allocation automatically. All memory allocation from within a dcmd must use UM_GC so that if the dcmd is interrupted by the user, the debugger can garbage-collect the memory.

mdb_zalloc is like mdb_alloc, but the allocated memory is filled with zeroes before returning it to the caller. No guarantees are made about the initial contents of memory returned by mdb_alloc.mdb_free is used to free previously allocated memory (unless it was allocated UM_GC). The buffer address and size must exactly match the original allocation. It is not legal to free only a portion of an allocation with mdb_free. It is not legal to free an allocation more than once. An allocation of zero bytes always returns NULL; freeing a NULL pointer with size zero always succeeds.

mdb printf

```
void mdb printf(const char *format, ...);
```

Print formatted output using the specified format string and arguments. Module writers should use mdb_printf for all output, except for warning and error messages. This function automatically triggers the built-in output pager when appropriate. The mdb_printf function is similar to printf(3C), with certain exceptions: the %C, %S, and %ws specifiers for wide character strings are not supported, the %f floating-point format is not supported, the %e, %E, %g, and %G specifiers for alternative double formats produce only a single style of output, and precision specifications of the form %.n are not supported. The list of specifiers that are supported follows:

Flag Specifiers

- \%\# If the \# sign is found in the format string, this selects the alternate form of the given format. Not all formats have an alternate form; the alternate form is different depending on the format. Refer to the format descriptions below for details on the alternate format.
- %+ When printing signed values, always display the sign (prefix with either '+' or '-'). Without %+, positive values have no sign prefix, and negative values have a '-' prefix prepended to them.
- %- Left-justify the output within the specified field width. If the width of the output is less than the specified field width, the output will be padded with blanks on the right-hand side. Without %-, values are right-justified by default.
- Zero-fill the output field if the output is right-justified and the width of the output is less than the specified field width. Without %0, right-justified values are prepended with blanks in order to fill the field.

Field Width Specifiers

- %n Field width is set to the specified decimal value.
- \$? Field width is set to the maximum width of a hexadecimal pointer value. This is 8 in an ILP32 environment, and 16 in an LP64 environment.
- %* Field width is set to the value specified at the current position in the argument list. This value is assumed to be an int. Note that in the 64-bit compilation environment, it may be necessary to cast long values to int.

Integer Specifiers

- %h Integer value to be printed is a short.
- **%1** Integer value to be printed is a long.
- \$11 Integer value to be printed is a long long.

Terminal Attribute Specifiers

If standard output for the debugger is a terminal, and terminal attributes can be obtained by the terminfo database, the following terminal escape constructs can be used:

- %<n>Enable the terminal attribute corresponding to n. Only a single attribute can be enabled with each instance of %<>.
- %</ Disable the terminal attribute corresponding to n. Note that in the case of reverse video, dim text, and bold text, the terminal codes to disable these attributes might be identical. Therefore, it might not be possible to disable these attributes independently of one another.</p>

If no terminal information is available, each terminal attribute construct is ignored by mdb_printf. For more information on terminal attributes, see terminfo(4). The available terminfo attributes are:

- Alternate character set
- b Bold text
- d Dim text
- r Reverse video
- s Best standout capability
- u Underlining

Format Specifiers

- %% The '%' symbol is printed.
- *a Prints an address in symbolic form. The minimum size of the value associated with *a is a uintptr_t; specifying *la is not necessary. If address-to-symbol conversion is on, the debugger will attempt to convert the address to a symbol name followed by an offset in the current output radix and print this string; otherwise, the value is printed in the default output radix. If *#a is used, the alternate format adds a ':' suffix to the output.
- This format is identical to %a, except when an address cannot be converted to a symbol name plus an offset, nothing is printed. If %#A is used, the alternate format prints a '?' when address conversion fails.
- Decode and print a bit field in symbolic form. This specifier expects two consecutive arguments: the bit field value (int for %b, long for %lb, and so forth), and a pointer to an array of mdb_bitmask_t structures:

```
} mdb_bitmask_t;
```

The array should be terminated by a structure whose bm_name field is set to NULL. When %b is used, the debugger reads the value argument, then iterates through each mdb_bitmask structure checking to see if:

```
(value & bitmask->bm_mask) == bitmask->bm_bits
```

If this expression is true, the bm_name string is printed. Each string printed is separated by a comma. The following example shows how %b can be used to decode the t_flag field in a kthread_t:

```
const mdb_bitmask_t t_flag_bits[] = {
    { "T_INTR_THREAD", T_INTR_THREAD, T_INTR_THREAD },
    { "T_WAKEABLE", T_WAKEABLE, T_WAKEABLE },
    { "T_TOMASK", T_TOMASK, T_TOMASK },
    { "T_TALLOCSTK", T_TALLOCSTK, T_TALLOCSTK },
    ...
    { NULL, 0, 0 }
};

void
thr_dump(kthread_t *t)
{
    mdb_printf("t_flag = <%hb>\n", t->t_flag, t_flag_bits);
    ...
}
```

If t_flag was set to 0x000a, the function would print:

```
t_flag = <T_WAKEABLE,T_TALLOCSTK>
```

If %#b is specified, the union of all bits that were not matched by an element in the bitmask array is printed as a hexadecimal value following the decoded names.

- %c Print the specified integer as an ASCII character.
- Print the specified integer as a signed decimal value. Same as %i. If %#d is specified, the alternate format prefixes the value with '0t'.
- *e Print the specified double in the floating-point format [+/-]d.dddddde[+/-]dd, where there is one digit before the radix character, seven digits of precision, and at least two digits following the exponent.
- Frint the specified double using the same rules as %e, except that the exponent character will be 'E' instead of 'e'.
- Print the specified double in the same floating-point format as %e, but with sixteen digits of precision. If %11g is specified, the argument is expected to be of type long double (quad-precision floating-point value).
- %G Print the specified double using the same rules as %g, except that the exponent character will be 'E' instead of 'e'.
- %i Print the specified integer as a signed decimal value. Same as %d. If %#i is specified, the alternate format prefixes the value with '0t'.

- Print the specified 32-bit unsigned integer as an Internet IPv4 address in dotted-decimal format (for example, the hexadecimal value 0xffffffff would print as 255.255.255.255).
- %m Print a margin of whitespace. If no field is specified, the default output margin width is used; otherwise, the field width determines the number of characters of white space that are printed.
- Print the specified integer as an unsigned octal value. If \$#0 is used, the alternate format prefixes the output with '0'.
- %p Print the specified pointer (void *) as a hexadecimal value.
- Print the specified integer as a signed octal value. If \$#0 is used, the alternate format prefixes the output with '0'.
- Frint the specified integer as an unsigned value in the current output radix. The user can change the output radix using the \$d dcmd. If \$\#\r is specified, the alternate format prefixes the value with the appropriate base prefix: '0i' for binary, '0o' for octal, '0t' for decimal, or '0x' for hexadecimal.
- R Print the specified integer as a signed value in the current output radix. If %#R is specified, the alternate format prefixes the value with the appropriate base prefix.
- %s Print the specified string (char *). If the string pointer is NULL, the string '<NULL>' is printed.
- Advance one or more tab stops. If no width is specified, output advances to the next tab stop; otherwise the field width determines how many tab stops are advanced.
- *T Advance the output column to the next multiple of the field width. If no field width is specified, no action is taken. If the current output column is not a multiple of the field width, white space is added to advance the output column.
- %u Print the specified integer as an unsigned decimal value. If %#u is specified, the alternate format prefixes the value with '0t'.
- Print the specified integer as a hexadecimal value. The characters a-f are used as the digits for the values 10-15. If \$\#\times\$ is specified, the alternate format prefixes the value with '0x'.
- Print the specified integer as a hexadecimal value. The characters A-F are used as the digits for the values 10-15. If \$#X is specified, the alternate format prefixes the value with '0X'.
- %Y The specified time_t is printed as the string 'year month day HH:MM:SS'.

mdb_snprintf

```
size_t mdb_snprintf(char *buf, size_t len, const char *format, ...);
```

Construct a formatted string based on the specified format string and arguments, and store the resulting string into the specified <code>buf</code>. The mdb_snprintf function accepts the same format specifiers and arguments as the mdb_printf function. The <code>len</code> parameter specifies the size of <code>buf</code> in bytes. No more than <code>len</code> - 1 formatted bytes are placed in <code>buf</code>; mdb_snprintf always terminates <code>buf</code> with a null byte. The function returns the number of bytes required for the complete formatted string, not including the terminating null byte. If the <code>buf</code> parameter is NULL and <code>len</code> is set to zero, the function will not store any characters to <code>buf</code> and returns the number of bytes required for the complete formatted string; this technique can be used to determine the appropriate size of a buffer for dynamic memory allocation.

mdb warn

```
void mdb_warn(const char *format, ...);
```

Print an error or warning message to standard error. The mdb_warn function accepts a format string and variable argument list that can contain any of the specifiers documented for mdb_printf. However, the output of mdb_warn is sent to standard error, which is not buffered and is not sent through the output pager or processed as part of a dcmd pipeline. All error messages are automatically prefixed with the string "mdb:".

In addition, if the format parameter does not contain a newline (\n) character, the format string is implicitly suffixed with the string ": %s\n", where %s is replaced by the error message string corresponding to the last error recorded by a module API function. For example, the following source code:

mdb flush

```
void mdb_flush(void);
```

Flush all currently buffered output. Normally, mdb's standard output is line-buffered; output generated using mdb_printf is not flushed to the terminal (or other standard output destination) until a newline is encountered, or at the end of the current dcmd. However, in some situations you might want to explicitly flush standard output prior to printing a newline; mdb_flush can be used for this purpose.

mdb_nhconvert

```
void mdb_nhconvert(void *dst, const void *src, size_t nbytes);
```

Convert a sequence of nbytes bytes stored at the address specified by src from network byte order to host byte order and store the result at the address specified by dst. The src and dst parameters may be the same, in which case the object is converted in place. This function may be used to convert from host order to network order or from network order to host order, since the conversion is the same in either case.

mdb_dumpptr and mdb_dump64

These functions can be used to generate formatted hexadecimal and ASCII data dumps that are printed to standard output. Each function accepts an addr parameter specifying the starting location, a nbytes parameter specifying the number of bytes to display, a set of flags described below, a func callback function to use to read the data to display, and a data parameter that is passed to each invocation of the callback func as its last argument. The functions are identical in every regard except that mdb_dumpptr uses uintptr_t for its address parameters and mdb_dump64 uses uint64_t. This distinction is useful when combining mdb_dump64 with mdb_pread, for example. The built-in ::dump dcmd uses these functions to perform its data display.

The flags parameter should be the bitwise OR of one or more of the following values:

MDB_DUMP_RELATIVE Number lines relative to the start address instead of with the explicit address of each line.

MDB_DUMP_ALIGN	Align the output at a paragraph boundary.
MDB_DUMP_PEDANT	Display full-width addresses instead of truncating the address to fit the output in 80 columns.
MDB_DUMP_ASCII	Display ASCII values next to the hexadecimal data.
MDB_DUMP_HEADER	Display a header line about the data.
MDB_DUMP_TRIM	Only read from and display the contents of the specified addresses, instead of reading and printing entire lines.
MDB_DUMP_SQUISH	Elide repeated lines by placing a "*" on a line that is a repeat of the previous line.
MDB_DUMP_NEWDOT	Update the value of dot to the address beyond the last address read by the function.
MDB_DUMP_ENDIAN	Adjust for endian-ness. This option assumes that the word size is equal to the current group size, specified by MDB_DUMP_GROUP. This option will always turn off alignment, headers, and ASCII display to avoid confusing output. If MDB_DUMP_TRIM is set with MDB_DUMP_ENDIAN, the number of bytes dumped will be rounded down to the nearest word size bytes.
MDB_DUMP_WIDTH(width)	Increase the number of 16-byte paragraphs per line that are displayed. The default value of <i>width</i> is one, and the maximum value is 16.
MDB_DUMP_GROUP(group)	Set the byte group size to <i>group</i> . The default <i>group</i> size is four bytes. The <i>group</i> size must be a power of two that divides the line width.

mdb one bit

```
const char *mdb_one_bit(int width, int bit, int on);
```

The mdb_one_bit function can be used to print a graphical representation of a bit field in which a single bit of interest is turned on or off. This function is useful for creating verbose displays of bit fields similar to the output from snoop(1M) -v. For example, the following source code:

Each bit in the bit field is printed as a period (.), with each 4-bit sequence separated by a white space. The bit of interest is printed as 1 or 0, depending on the setting of the on parameter. The total width of the bit field in bits is specified by the width parameter, and the bit position of the bit of interest is specified by the bit parameter. Bits are numbered starting from zero. The function returns a pointer to

an appropriately sized, null-terminated string containing the formatted bit representation. The string is automatically garbage-collected upon completion of the current dcmd.

mdb_inval_bits

```
const char *mdb_inval_bits(int width, int start, int stop);
```

The mdb_inval_bits function is used, along with mdb_one_bit, to print a graphical representation of a bit field. This function marks a sequence of bits as invalid or reserved by displaying an 'x' at the appropriate bit location. Each bit in the bit field is represented as a period (.), except for those bits in the range of bit positions specified by the start and stop parameters. Bits are numbered starting from zero. For example, the following source code:

```
mdb_printf("%s = reserved\n", mdb_inval_bits(8, 7, 7));
produces this output:
x... = reserved
```

The function returns a pointer to an appropriately sized, null-terminated string containing the formatted bit representation. The string is automatically garbage-collected upon completion of the current dcmd.

mdb_inc_indent and mdb_dec_indent

```
ulong_t mdb_inc_indent(ulong_t n);
ulong_t mdb_dec_indent(ulong_t n);
```

These functions increment and decrement the numbers of columns that MDB will auto-indent with white space before printing a line of output. The size of the delta is specified by n, a number of columns. Each function returns the previous absolute value of the indent. Attempts to decrement the indent below zero have no effect. Following a call to either function, subsequent calls to mdb_printf are indented appropriately. If the dcmd completes or is forcibly terminated by the user, the indent is restored automatically to its default setting by the debugger.

mdb_eval

```
int mdb_eval(const char *s);
```

Evaluate and execute the specified command string s, as if it had been read from standard input by the debugger. This function returns 0 for success, or -1 for error. mdb_eval fails if the command string contains a syntax error, or if the command string executed by mdb_eval is forcibly aborted by the user using the pager or by issuing an interrupt.

mdb_set_dot and mdb_get_dot

```
void mdb_set_dot(uintmax_t dot);
uintmax_t mdb_get_dot(void);
```

Set or get the current value of dot (the "." variable). Module developers might want to reposition dot so that, for example, it refers to the address following the last address read by the dcmd.

mdb_get_pipe

```
void mdb_get_pipe(mdb_pipe_t *p);
```

Retrieve the contents of the pipeline input buffer for the current dcmd. The mdb_get_pipe function is intended to be used by dcmds that want to consume the complete set of pipe input and execute only once, instead of being invoked repeatedly by the debugger for each pipe input element. Once mdb_get_pipe is invoked, the dcmd will not be invoked again by the debugger as part of the current command. This can be used, for example, to construct a dcmd that sorts a set of input values.

The pipe contents are placed in an array that is garbage-collected upon termination of the dcmd, and the array pointer is stored in p->pipe_data. The length of the array is placed in p->pipe_len. If the dcmd was not executed on the right-hand side of a pipeline (that is, the DCMD_PIPE flag was not set in its flags parameter), p->pipe_data is set to NULL and p->pipe_len is set to zero.

mdb_set_pipe

```
void mdb_set_pipe(const mdb_pipe_t *p);
```

Set the pipeline output buffer to the contents described by the pipe structure p. The pipe values are placed in the array p->pipe_data, and the length of the array is stored in p->pipe_len. The debugger makes its own copy of this information, so the caller must remember to free p->pipe_data if necessary. If the pipeline output buffer was previously non-empty, its contents are replaced by the new array. If the dcmd was not executed on the left side of a pipeline (that is, the DCMD_PIPE_OUT flag was not set in its flags parameter), this function has no effect.

mdb_get_xdata

```
ssize_t mdb_get_xdata(const char *name, void *buf, size_t nbytes);
```

Read the contents of the target external data buffer specified by name into the buffer specified by <code>buf</code>. The size of <code>buf</code> is specified by the <code>nbytes</code> parameter; no more than <code>nbytes</code> will be copied to the caller's buffer. The total number of bytes read will be returned upon success; -1 will be returned upon error. If the caller wants to determine the size of a particular named buffer, <code>buf</code> should be specified as NULL and <code>nbytes</code> should be specified as zero. In this case, <code>mdb_get_xdata</code> will return the total size of the buffer in bytes but no data will be read. External data buffers provide module writers access to target data that is not otherwise accessible through the module API. The set of named buffers exported by the current target can be viewed using the <code>::xdata</code> built-in dcmd.

Additional Functions

Additionally, module writers can use the following string(3C) and bstring(3C) functions. They are guaranteed to have the same semantics as the functions described in the corresponding Solaris man page.

```
        strcat()
        strcpy()
        strncpy()

        strchr()
        strcmp()
        strcmp()

        strncmp()
        strcasecmp()
        strncasecmp()

        strlen()
        bcmp()
        bcopy()

        bzero()
        bsearch()
        qsort()
```

Appendix A. Options

This appendix provides a reference for MDB command-line options.

Summary of Command-line Options

```
mdb [ -fkmuwyAFMS ] [ +o option ] [ -p pid ] [ -s distance]
        [ -I path ] [ -L path ] [ -P prompt ] [ -R root ]
        [ -V dis-version ] [ object [ core ] | core | suffix ]
```

The following options are supported:

- A Disables automatic loading of **mdb** modules. By default, **mdb** attempts to load debugger modules corresponding to the active shared libraries in a user process or core file, or to the loaded kernel modules in the live operating system or an operating system crash dump.
- Forcibly takes over the specified user process, if necessary. By default, **mdb** refuses to attach to a user process that is already under the control of another debugging tool, such as truss(1). With the F option, **mdb** attaches to these processes anyway. This can produce unexpected interactions between **mdb** and the other tools attempting to control the process.
- Force raw file debugging mode. By default, **mdb** attempts to infer whether the object and core file operands refer to a user executable and core dump or to a pair of operating system crash dump files. If the file type cannot be inferred, the debugger will default to examining the files as plain binary data. The foption forces **mdb** to interpret the arguments as a set of raw files to examine
- I Sets default path for locating macro files. Macro files are read using the \$< or \$<< dcmds. The path is a sequence of directory names delimited by colon (:) characters. The I include path and L library path (see below) can also contain any of the following tokens:
 - %i Expands to the current instruction set architecture (ISA) name: sparc, sparcv9, i386, or amd64.
 - %o Expands to the old value of the path being modified. This is useful for appending or prepending directories to an existing path.
 - %p Expands to the current platform string (either **uname** i or the platform string stored in the process core file or crash dump).
 - %r Expands to the path name of the root directory. An alternate root directory can be specified using the R option. If no R option is present, the root directory is derived dynamically from the path to the **mdb** executable itself. For example, if /bin/mdb is executed, the root directory is /. If /net/hostname/bin/mdb were executed, the root directory would be derived as /net/hostname.
 - %t Expands to the name of the current target. This is either the literal string 'proc' (a user process or user process core file), or 'kvm' (a kernel crash dump or the live operating system).

The default include path for 32-bit **mdb** is:%r/usr/platform/%p/lib/adb:%r/usr/lib/adb

The default include path for 64-bit **mdb** is:%r/usr/platform/%p/lib/adb/%i:%r/usr/lib/adb/%i

k Forces kernel debugging mode. By default, **mdb** attempts to infer whether the object and core file operands refer to a user executable and core dump, or to a pair of operating system crash dump

- files. The k option forces **mdb** to assume these files are operating system crash dump files. If no object or core operand is specified, but the k option is specified, **mdb** defaults to an object file of /dev/ksyms and a core file of /dev/kmem. Access to /dev/kmem is restricted to group sys.
- K Load **kmdb**, stop the live running operating system kernel, and proceed to the **kmdb** debugger prompt. This option should only be used on the system console, as the subsequent **kmdb** prompt will appear on the system console.
- L Sets default path for locating debugger modules. Modules are loaded automatically on startup or by using the ::load dcmd. The path is a sequence of directory names delimited by colon (:) characters. The L library path can also contain any of the tokens shown for I above.
- m Disables demand-loading of kernel module symbols. By default, mdb processes the list of loaded kernel modules and performs demand loading of per-module symbol tables. If the m option is specified, mdb does not attempt to process the kernel module list or provide per-module symbol tables. As a result, mdb modules corresponding to active kernel modules are not loaded on startup.
- M Preloads all kernel module symbols. By default, **mdb** performs demand-loading for kernel module symbols: the complete symbol table for a module is read when an address is that module's text or data section is referenced. With the M option, **mdb** loads the complete symbol table of all kernel modules during startup.
- o optEmables the specified debugger option. If the o form of the option is used, the specified option is disabled. Unless noted below, each option is off by default. **mdb** recognizes the following option arguments:
 - adb Enable stricter adb(1) compatibility. The prompt is set to the empty string and many **mdb** features, such as the output pager, are disabled.
 - array_Seetth_dideriaultilimit on the number of array members that ::print will display. If limit is the special token none, all array members will be displayed by default.
 - array_Stt_thmidefairhtilimit on the number of characters that ::print will attempt to display as an ASCII string when printing a char array. If limit is the special token none, the entire char array will be displayed as a string by default.
 - follow Setxthe_dubdeggenote havior for following an exec(2) system call. The mode should be one of the following named constants:

ask	If stdout is a terminal device, the debugger will stop after the exec(2) system call has returned and then prompt the user to decide whether to follow the exec or stop. If stdout is not a terminal device, the ask mode will default to stop.
follow	The debugger will follow the exec by automatically continuing the target process and resetting all of its mappings and symbol tables based on the new executable. The follow behavior is discussed in more detail under Interaction With exec.
stop	The debugger will stop following return from the exec system call. The stop behavior is discussed in more detail under Interaction With exec.

follow Setotke debutegenoblehavior for following a fork(2), fork1(2), or vfork(2) system call. The mode should be one of the following named constants:

ask	If stdout is a terminal device, the debugger will stop after the fork system call
	has returned and then prompt the user to decide whether to follow the parent or
	child. If stdout is not a terminal device, the ask mode will default to parent.

parent	The debugger will follow the parent process, and will detach from the child process and set it running.
	The debugger will follow the child process, and will detach from the parent process and set it running.

ignore Ehf debugger does not exit when an EOF sequence (^D) is entered at the terminal. The ::quit dcmd must be used to quit.

nostopDo not stop a user process when attaching to it when the -p option is specified or when the ::attach or :A dcmds are applied. The nostop behavior is described in more detail under Process Attach and Release.

pager The output pager is enabled (default).

repeatlast NEWLINE is entered as the complete command at the terminal, **mdb** repeats the previous command with the current value of dot. This option is implied by o **adb**.

showlMiDB provides support for symbol naming and identification in user applications that make use of link maps other than LM_ID_BASE and LM_ID_LDSO, as described in Symbol Name Resolution. Symbols on link maps other than LM_ID_BASE or LM_ID_LDSO will be shown as LMlmid`library`symbol, where lmid is the link-map ID in the default output radix (16). The user may optionally configure MDB to show the link-map ID scope of all symbols and objects, including those associated with LM_ID_BASE and LM_ID_LDSO, by enabling the showlmid option. Built-in dcmds that deal with object file names will display link-map IDs according to the value of showlmid above, including ::nm, ::mappings, \$m, and ::objects.

- p pidAttaches to and stops the specified process id. **mdb** uses the /proc/pid/object/a.out file as the executable file path name.
- P Sets the command prompt. The default prompt is '> '.
- R Sets root directory for path name expansion. By default, the root directory is derived from the path name of the **mdb** executable itself. The root directory is substituted in place of the %r token during path name expansion.
- s distance. By default, **mdb** sets the distance to zero, which enables a smart-matching mode. Each ELF symbol table entry includes a value V and size S, representing the size of the function or data object in bytes. In smart mode, **mdb** matches an address A with the given symbol if A is in the range [V, V + S]. If any non-zero distance is specified, the same algorithm is used, but S in the given expression is always the specified absolute distance and the symbol size is ignored.
- S Suppresses processing of the user's ~/.mdbrc file. By default, **mdb** reads and processes the macro file.mdbrc if one is present in the user's home directory, as defined by \$HOME. If the S option is present, this file is not read.
- Forces user debugging mode. By default, **mdb** attempts to infer whether the object and core file operands refer to a user executable and core dump, or to a pair of operating system crash dump files. The u option forces **mdb** to assume these files are not operating system crash dump files.
- Unload **kmdb** if it is loaded. You should unload **kmdb** when it is not in use to release the memory used by the kernel debugger back to the free memory available to the operating system.

- V Sets disassembler version. By default, **mdb** attempts to infer the appropriate disassembler version for the debug target. The disassembler can be set explicitly using the V option. The ::disasms dcmd lists the available disassembler versions.
- w Opens the specified object and core files for writing.
- y Sends explicit terminal initialization sequences for tty mode. Some terminals require explicit initialization sequences to switch into a tty mode. Without this initialization sequence, terminal features such as standout mode might not be available to **mdb**.

Operands

The following operands are supported:

object Specifies an ELF format object file to examine. **mdb** provides the ability to examine and edit ELF format executables (ET_EXEC), ELF dynamic library files (ET_DYN), ELF relocatable object files (ET_REL), and operating system unix.X symbol table files.

Specifies an ELF process core file (ET_CORE), or an operating system crash dump vmcore.X file. If an ELF core file operand is provided without a corresponding object file, **mdb** will attempt to infer the name of the executable file that produced the core using several different algorithms. If no executable is found, **mdb** will still execute, but some symbol information may be unavailable.

Specifies the numerical suffix representing a pair of operating system crash dump files. For example, if the suffix is '3', **mdb** infers that it should examine the files 'unix.3' and 'vmcore.3'. The string of digits will not be interpreted as a suffix if an actual file of the same name is present in the current directory.

Exit Status

The following exit values are returned:

- 0 Debugger completed execution successfully.
- 1 A fatal error occurred.
- 2 Invalid command line options were specified.

Environment Variables

The following environment variables are supported:

HISTSIZE This variable is used to determine the maximum length of the command history list. If this variable is not present, the default length is 128.

HOME This variable is used to determine the pathname of the user's home directory, where a .mdbrc file may reside. If this variable is not present, no .mdbrc processing will occur.

SHELL This variable is used to determine the pathname of the shell used to process shell escapes requested using the! meta-character. If this variable is not present, /bin/sh is used.

Appendix B. Notes

Warnings

The following warning information applies to the use of MDB.

Use of the Error Recovery Mechanism

The debugger and its dmods execute in the same address space, and thus it is quite possible that a buggy dmod can cause MDB to dump core or otherwise misbehave. The MDB **resume** capability, described in Signal Handling, provides a limited recovery mechanism for these situations. However, it is not possible for MDB to know definitively whether the dmod in question has corrupted only its own state, or the debugger's global state. Therefore a **resume** operation cannot be guaranteed to be safe, or to prevent a subsequent crash of the debugger. The safest course of action following a **resume** is to save any important debug information, and then quit and restart the debugger.

Use of the Debugger to Modify the Live Operating System

The use of the debugger to modify (that is, write to) the address space of live running operating system is extremely dangerous, and may result in a system panic in the event the user damages a kernel data structure.

Use of kmdb to Stop the Live Operating System

The use of **kmdb** to stop the live operating system using mdb –K or by setting a breakpoint in the live operating system is intended for use by developers and not on production systems. When the operating system kernel is stopped by **kmdb**, operating system services and networking are not executing, and other systems on the network that depend upon the target system will not be able to contact the target system.

Notes

Limitations on Examining Process Core Files

MDB does not provide support for examining process core files that were generated by a release of the Solaris operating system preceding Solaris 2.6. If a core file from one operating system release is examined on a different operating system release, the run-time link-editor debugging interface (librtld_db) may not be able to initialize. In this case, symbol information for shared libraries will not be available. Furthermore, since shared mappings are not present in user core files, the text section and read-only data of shared libraries may not match the data that was present in the process at the time it dumped core. Core files from Solaris x86 systems may not be examined on Solaris SPARC systems, and vice-versa.

Limitations on Examining Crash Dump Files

Crash dumps from Solaris 7 and earlier releases may only be examined with the aid of the libkvm from the corresponding operating system release. If a crash dump from one operating system release is examined using the dmods from a different operating system release, changes in the kernel implementation may prevent some dcmds or walkers from working properly. MDB will issue a warning message if it detects this condition. Crash dumps from Solaris x86 systems may not be examined on Solaris SPARC systems, and vice-versa.

Relationship Between 32-bit and 64-bit Debugger

MDB provides support for debugging both 32-bit and 64-bit programs. Once it has examined the target and determined its data model, MDB will automatically re-execute the mdb binary that has the same data model as the target, if necessary. This approach simplifies the task of writing debugger modules, because the modules that are loaded will use the same data model as the primary target. Only the 64-bit debugger may be used to debug 64-bit target programs. The 64-bit debugger can only be used on a system that is running the 64-bit operating environment.

Limitations on Memory Available to kmdb

The memory available to **kmdb** is allocated when the debugger is loaded, and cannot be expanded after that point in time. If debugger commands attempt to allocate more memory than is available, they will not be able to execute. The debugger will attempt to gracefully recover from low memory situations, but may be forced to terminate the system under dire circumstances. System memory constraints are especially acute on x86 platforms that use the 32-bit operating system kernel.

Developer Information

The mdb(1) man page provides a detailed description of built-in **mdb** features for easy developer reference. The header file <sys/mdb_modapi.h> contains prototypes for the functions in the MDB Module API, and the SUNWmdbdm package provides source code for an example module in the directory /usr/demo/mdb.

Appendix C. Transition From adb and kadb

The transition from using the legacy adb(1) utility to using mdb(1) is relatively simple: MDB provides evolutionary compatibility for the adb syntax, built-in commands, and command-line options. MDB attempts to provide compatibility for all existing adb(1) features, but it is not bug-for-bug compatible with adb(1). This appendix briefly discusses several features of adb(1) that are not precisely emulated by mdb(1) in order to guide users to the new functionality

Command-line Options

MDB provides a superset of the command-line options recognized by adb(1). All the adb(1) options are supported and have the same meaning as before. The /usr/bin/adb pathname is delivered as a link that invokes mdb(1), and automatically enables enhanced adb(1) compatibility mode. Executing the /usr/bin/adb link is equivalent to executing mdb with the o adb option, or executing ::set -o adb once the debugger has started.

Syntax

The MDB language adheres to the same syntax as the adb(1) language, in order to provide compatibility for legacy macros and script files. New MDB dcmds use the extended form ::name, in order to distinguish them from legacy commands that are prefixed with either : or \$. Expressions can also be evaluated on the right-hand side of a dcmd name by enclosing them in square brackets preceded by a dollar sign (\$[]). Similar to adb(1), an input line that begins with an exclamation mark (!) indicates that the command line should be executed by the user's shell. In MDB, a debugger command may also be suffixed with an exclamation mark to indicate that its output should be piped to the shell command following the exclamation mark.

In adb(1), binary operators are left associative and have lower precedence than unary operators. Binary operators are evaluated in strict left-to-right order on the input line. In MDB, binary operators are left associative and have lower precedence than unary operators, but the binary operators operate in order of precedence according to the table in Binary Operators. The operators conform to the order of precedence in ANSI C. Legacy adb(1) macro files that do not explicitly parenthesize ambiguous expressions may need to be updated to work with MDB. For example, in adb the following command evaluates to the integer value nine:

```
$ echo "4-1*3=X" | adb
```

In MDB, as in ANSI C, operator "*" has higher precedence than "-" and therefore the result is the integer value one:

Watchpoint Length Specifier

The watchpoint length specifier syntax recognized by MDB is different from the syntax described in adb(1). In particular, the adb watchpoint commands : w, :a, and :p allow an integer length in bytes to be

inserted between the colon and the command character. In MDB, the count should be specified following the initial address as a repeat count. Stated simply, these adb(1) commands:

```
123:456w
123:456a
123:456p
are specified in MDB as
123,456:w
123,456:a
```

123,456:p

The MDB: :wp dcmd provides more complete facilities for creating user process watchpoints. Similarly, the legacy **kadb** length modifier command \$1 is not supported. Therefore, the watchpoint size should be specified to each: :wp command used in **kmdb**.

Address Map Modifier

The adb(1) commands to modify segments of the virtual address map and object file map are not present in MDB. Specifically, the /m, /*m, ?m, and ?*m format specifiers are not recognized or supported by MDB. These specifiers were used to manually modify the valid addressable range of the current object and core files. MDB properly recognizes the addressable range of such files automatically, and updates the ranges when a live process is being debugged, so these commands are no longer necessary.

Output

The precise text output form of some commands is different in MDB. Macro files are formatted using the same basic rules, but shell scripts that depend on the precise character-by-character output of certain commands may need to change. Users who have shell scripts that parse the output of adb commands will need to revalidate and update such scripts as part of the transition to MDB.

Deferred Breakpoints

The legacy **kadb** utility supported a syntax for deferred breakpoints that was incompatible with the existing **adb** syntax. These deferred breakpoints were specified using the syntax *module#symbol*:b in **kadb**. To set a deferred breakpoint in **kmdb**, use the MDB ::bp dcmd as described in Chapter 6, Execution Control.

I/O Port Access

The legacy **kadb** utility provided access to I/O ports on x86 systems using the :i and :o commands. These commands are not supported in **mdb** or **kmdb**. Access to I/O ports on x86 systems is provided by the ::in and ::out commands.

Appendix D. Transition From crash

The transition from using the legacy **crash** utility to using mdb(1) is relatively simple: MDB provides most of the "canned" crash commands. The additional extensibility and interactive features of MDB allow the programmer to explore aspects of the system not examined by the current set of commands. This appendix briefly discusses several features of **crash** and provides pointers to equivalent MDB functionality.

Command-line Options

The crash d, n, and w command-line options are not supported by **mdb**. The crash dump file and name list (symbol table file) are specified as arguments to **mdb** in the order of name list, crash dump file. To examine the live kernel, the mdb k option should be specified with no additional arguments. Users who want to redirect the output of mdb to a file or other output destination, should either employ the appropriate shell redirection operator following the **mdb** invocation on the command line, or use the ::log built-in dcmd.

Input in MDB

In general, input in MDB is similar to **crash**, except that function names (in MDB, dcmd names) are prefixed with "::". Some MDB dcmds accept a leading expression argument that precedes the dcmd name. Like **crash**, string options can follow the dcmd name. If a! character follows a function invocation, MDB will also create a pipeline to the specified shell pipeline. All immediate values specified in MDB are interpreted in hexadecimal by default. The radix specifiers for immediate values are different in **crash** and MDB as shown in Table D-1:

Table D.1. Radix Specifiers

crash	mdb	Radix
0x	0x	hexadecimal (base 16)
0d	Ot	decimal (base 10)
0b	0i	binary (base 2)

Many **crash** commands accepted slot numbers or slot ranges as input arguments. The Solaris operating system is no longer structured in terms of slots, so MDB dcmds do not provide support for slot-number processing.

Functions

crash function	mdb dcmd	Comments
?	::dcmds	List available functions.
!command	!command	Escape to the shell and execute command.
base	=	In mdb, the = format character can be used to convert the left-hand expression value to any of the known formats. Formats for octal, decimal, and hexadecimal are provided.
callout	::callout	Print the callout table.
class	::class	Print scheduling classes.
cpu	::cpuinfo	

crash function	mdb dcmd	Comments
		Print information about the threads dispatched on the system CPUs. If the contents of a particular CPU structure are needed, the user should apply the \$ <cpu address="" cpu="" in="" macro="" mdb.<="" td="" the="" to=""></cpu>
help	::help	Print a description of the named dcmd, or general help information.
kfp	::regs	The mdb ::regs dcmd displays the complete kernel register set, including the current stack frame pointer. The \$C dcmd can be used to display a stack backtrace including frame pointers.
kmalog	::kmalog	Display events in kernel memory allocator transaction log.
kmastat	::kmastat	Print kernel memory allocator transaction log.
kmausers	::kmausers	Print information about the medium and large users of the kernel memory allocator that have current memory allocations.
mount	::fsinfo	Print information about mounted file systems.
nm	::nm	Print symbol type and value information.
od	::dump	Print a formatted memory dump of a given region. In mdb, ::dump displays a mixed ASCII and hexadecimal display of the region.
proc	::ps	Print a table of the active processes.
quit	::quit	Quit the debugger.
rd	::dump	Print a formatted memory dump of a given region. In mdb, ::dump displays a mixed ASCII and hexadecimal display of the region.
redirect	::log	In mdb, output for input and output can be globally redirected to a log file using ::log.
search	::kgrep	In mdb, the ::kgrep dcmd can be used to search the kernel's address space for a particular value. The pattern match built-in dcmds can also be used to search the physical, virtual, or object files address spaces for patterns.
stack	::stack	The current stack trace can be obtained using ::stack. The stack trace of a particular kernel thread can be determined using the ::findstack dcmd. A memory dump of the current stack can be obtained using the / or ::dump dcmds and the current stack pointer. The \$ <stackregs a="" applied="" be="" can="" macro="" obtain="" per-frame="" pointer="" register="" saved="" stack="" td="" the="" to="" values.<=""></stackregs>
status	::status	Display status information about the system or dump being examined by the debugger.
stream	::stream	The mdb ::stream dcmd can be used to format and display the structure of a particular kernel STREAM. If the list of active STREAM structures is needed, the user should execute ::walk stream_head_cache in mdb and pipe the resulting addresses to an appropriate formatting dcmd or macro.
strstat	::kmastat	The ::kmastat dcmd displays a superset of the information reported by the strstat function.
trace	::stack	The current stack trace can be obtained using ::stack. The stack trace of a particular kernel thread can be determined using the ::findstack dcmd. A memory dump of the current stack can be obtained using the / or ::dump dcmds and the current stack pointer. The \$ <stackregs a="" applied="" be="" can="" macro="" obtain="" per-frame="" pointer="" register="" saved="" stack="" td="" the="" to="" values.<=""></stackregs>

Transition From crash

crash function	mdb dcmd	Comments
var	\$ <v< td=""><td>Print the tunable system parameters in the global var structure.</td></v<>	Print the tunable system parameters in the global var structure.
vfs	::fsinfo	Print information about mounted file systems.
vtop	::vtop	Print the physical address translation of the given virtual address.