

**The DiskSim Simulation Environment  
Version 3.0 Reference Manual**

John S. Bucy, Gregory R. Ganger, and Contributors  
January 2003  
CMU-CS-03-102

School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213

**Abstract**

*DiskSim is an efficient, accurate and highly-configurable disk system simulator developed to support research into various aspects of storage subsystem architecture. It includes modules that simulate disks, intermediate controllers, buses, device drivers, request schedulers, disk block caches, and disk array data organizations. In particular, the disk drive module simulates modern disk drives in great detail and has been carefully validated against several production disks (with accuracy that exceeds any previously reported simulator).*

*This manual describes how to configure and use DiskSim, which has been made publicly available with the hope of advancing the state-of-the-art in disk system performance evaluation in the research community. The manual also briefly describes DiskSim's internal structure and various validation results.*

**Keywords:** storage system, disk simulator, disk model

# Contents

<b>Contents</b>	<b>i</b>
<b>1 Introduction</b>	<b>1</b>
1.1 What DiskSim Does . . . . .	1
1.2 What DiskSim Does Not Do . . . . .	1
1.3 Limitations and Advantages of Version 3.0 . . . . .	2
1.3.1 Diskmodel . . . . .	2
1.3.2 Libparam . . . . .	2
1.4 Known Bugs . . . . .	2
1.5 Organization of Manual . . . . .	2
1.5.1 Contributors . . . . .	2
<b>2 Running DiskSim</b>	<b>3</b>
2.1 Parameter Overrides . . . . .	3
2.2 Example Command Line . . . . .	4
<b>3 The Parameter File</b>	<b>5</b>
3.1 Global Block . . . . .	5
3.2 Stats Block . . . . .	6
3.2.1 Bus Statistics . . . . .	6
3.2.2 Controller Statistics . . . . .	6
3.2.3 Device Statistics . . . . .	7
3.2.4 iodriver Statistics . . . . .	7
3.2.5 Process-flow Statistics . . . . .	7
3.3 iosim Block . . . . .	7
3.4 I/O Subsystem Component Specifications . . . . .	8
3.4.1 Device Drivers . . . . .	8
3.4.2 Buses . . . . .	9
3.4.3 Controllers . . . . .	10
3.4.4 Storage Devices . . . . .	10
3.4.5 Disks . . . . .	11
3.4.6 Simple Disks . . . . .	18
3.4.7 Queue/Scheduler Subcomponents . . . . .	19
3.4.8 Disk Block Cache Subcomponents . . . . .	20
3.4.9 Memory Caches . . . . .	22
3.4.10 Cache Devices . . . . .	23
3.5 Component Instantiation . . . . .	24
3.6 I/O Subsystem Interconnection Specifications . . . . .	25
3.7 Rotational Synchronization of Devices . . . . .	25
3.8 Disk Array Data Organizations . . . . .	25
3.9 Process-Flow Parameters . . . . .	28
<b>4 Input Workloads: Traces and Synthetic Workloads</b>	<b>29</b>
4.1 Traces . . . . .	29
4.1.1 Default Format . . . . .	29
4.1.2 Adding Support For New Trace Formats . . . . .	29
4.2 Synthetic Workloads . . . . .	30
4.2.1 Configuration . . . . .	30
4.3 Incorporating DiskSim Into System-Level Simulators . . . . .	32

<b>5</b>	<b>The Output File</b>	<b>34</b>
5.1	The statdefs File . . . . .	34
5.2	Simulation Results . . . . .	35
5.2.1	Process-flow Statistics . . . . .	35
5.2.2	Validation Trace Statistics . . . . .	37
5.2.3	HPL Trace Statistics . . . . .	37
5.2.4	System-level Logical Organization Statistics . . . . .	38
5.2.5	I/O Driver Statistics . . . . .	40
5.2.6	Disk Statistics . . . . .	41
5.2.7	Controller Statistics . . . . .	43
5.2.8	Bus Statistics . . . . .	44
<b>6</b>	<b>Validation</b>	<b>45</b>
<b>A</b>	<b>Copyright notices for DiskSim</b>	<b>48</b>
A.1	Version 3.0 Copyright Addendum . . . . .	48
A.2	Version 2.0 Copyright Addendum . . . . .	48
A.3	Original (Version 1.0) Copyright Statement . . . . .	48
<b>B</b>	<b>Diskmodel</b>	<b>50</b>
B.1	Introduction . . . . .	50
B.2	Types and Units . . . . .	50
B.2.1	Three Zero Angles . . . . .	50
B.2.2	Two Zero Sectors . . . . .	50
B.2.3	Example . . . . .	50
B.3	API . . . . .	51
B.3.1	Disk-wide Parameters . . . . .	51
B.3.2	Layout . . . . .	51
B.3.3	Mechanics . . . . .	53
B.4	Model Configuration . . . . .	56
B.4.1	dm_disk . . . . .	56
B.4.2	dm_layout_g1 . . . . .	56
B.4.3	dm_mech_g1 . . . . .	58
	<b>References</b>	<b>61</b>

# 1 Introduction

Because of trends in both computer technology advancement (e.g., CPU speeds vs. disk access times) and application areas (e.g., on-demand video, global information access), storage system performance is becoming an increasingly large determinant of overall system performance. As a result, the need to understand storage performance under a variety of workloads is growing. Disk drives, which are still the secondary storage medium of choice, continue to expand in capacity and reliability while decreasing in unit cost, price/capacity, and power requirements. Performance characteristics also continue to change due to maturing technologies as well as new advances in materials, sensors, and electronics. New storage subsystem architectures may be needed to better exploit current and future generations of disk devices. The DiskSim simulation environment was developed as a tool for two purposes: understanding storage performance and evaluating new architectures.

## 1.1 What DiskSim Does

DiskSim is an efficient, accurate, highly-configurable storage system simulator. It is written in C and requires no special system software. It includes modules for many secondary storage components of interest, including device drivers, buses, controllers, adapters and disk drives. Some of the major functions (e.g., request queuing/scheduling, disk block caching, disk array data organizations) that can be present in several different components (e.g., operating system software, intermediate controllers, disk drives) have been implemented as separate modules that are linked into components as desired. Some of the component modules are highly detailed (e.g., the disk module), and the individual components can be configured and interconnected in a variety of ways. DiskSim has been used in a variety of published studies (and several unpublished studies) to understand modern storage subsystem performance [3, 22], to understand how storage performance relates to overall system performance [5, 2, 1], and to evaluate new storage subsystem architectures [21].

DiskSim has been validated both as part of a more comprehensive system-level model [5, 1] and as a standalone subsystem [22, 24, 18]. In particular, the disk module, which is extremely detailed, has been carefully validated against five different disk drives from three different manufacturers. The accuracy demonstrated exceeds that of any other disk simulator known to the authors (e.g., see [16]).

DiskSim can be driven by externally-provided I/O request traces or internally-generated synthetic workloads. Several trace formats have been used and new ones can be easily added. The synthetic trace generation module is quite flexible, particularly in the request arrival model (which can mimic an open process, a closed process or something in between). DiskSim was originally part of a larger, system-level model [5, 2] that modeled each request's interactions with executing processes, but has been separated out for public dissemination.<sup>1</sup> As a result, it can be integrated into full system simulators (e.g., simulators like SimOS [14]) with little difficulty.

## 1.2 What DiskSim Does Not Do

DiskSim, by itself, simulates and reports on only the performance-related aspects of the storage subsystem. It does not model the behavior of the other computer system components or interactions between them and the storage subsystem.<sup>2</sup> Because storage subsystem performance metrics are not absolute indicators of overall system performance (e.g., see [1]), promising architectures should be evaluated in the context of a more comprehensive system model or a real system. In such cases, DiskSim becomes one component of the full model, just as a storage subsystem is one component of a full system.

DiskSim models the performance behavior of disk systems, but does not actually save or restore data for each request. If such functionality is desired (as, for example, when building a full system simulator like SimOS), it can easily be provided independently of DiskSim, which will still provide accurate timings for I/O activity. See [6] for an example of this in the form of storage subsystem emulation.

<sup>1</sup>The system-level model includes several portions of a proprietary operating system, allowing it to achieve a close match to the real system behavior [2] but also preventing it from being publicly released.

<sup>2</sup>Actually, a rudimentary system model was kept in place to support the internal synthetic generation module. However, it should not be viewed as representative of any real system's behavior.

### 1.3 Limitations and Advantages of Version 3.0

DiskSim 3.0 builds on DiskSim 2.0 in several ways: parts of DiskSim's disk model have been moved into a separate library (diskmodel), the library has been re-integrated into DiskSim, the parameter-file infrastructure has been completely rewritten, and most of it has been moved into another library (libparam). The checkpoint-restore code has been removed for this release as it was not well maintained and created debugging difficulties. Support has been added for using one storage device as a cache for another.

There are still a number of limitations on the shape of a storage system topology. See Section 3.6 for details.

#### 1.3.1 Diskmodel

The impetus for the diskmodel transition was the desire to use DiskSim's disk model with a number of other software projects. The original interface between DiskSim's disk model and the rest of DiskSim was idiosyncratic to DiskSim, making it difficult to reuse. As a result, DiskSim's layout and mechanical code was often replicated in other software in an ad hoc fashion that resulted in numerous incompatible copies that were difficult to keep in sync.

Diskmodel introduces a clean, functional interface to all of the mechanical and layout computations. Integration of the new interface required a substantial overhaul of DiskSim with the advantage that, now, the implementations of the diskmodel functions can change independently of DiskSim and that diskmodel has no knowledge about DiskSim's internals embedded in it.

#### 1.3.2 Libparam

Libparam also unifies DiskSim's parameter-input code such that the same parser (in libparam) can be shared by disksim and diskmodel. This also makes it easy for applications using diskmodel to input the necessary disk specifications without copying large portions of DiskSim's input code.

Libparam introduces a new grammar for configuring DiskSim that is easier to use. It tolerates reordered parameters, unlike DiskSim 2.0, and generally provides greater assistance in identifying bugs in inputs.

### 1.4 Known Bugs

Multi-CPU configurations in the process-flow simulation do not work correctly; only the first CPU will have processes scheduled on it while the rest idle forever.

### 1.5 Organization of Manual

This manual describes how to configure and use DiskSim and how to interpret the output generated. Section 2 explains the DiskSim command line and how to execute the simulator given the appropriate configuration files. Section 3 describes how to describe a storage system to DiskSim. Section 4 describes how to provide an input workload of I/O requests to DiskSim – options include external traces, internally-generated synthetic workloads, and interactions with a larger simulation environment. Section 5 describes the contents of the output file. Section 6 provides validation data (all of which has been published previously [22, 24, 1]) showing that DiskSim accurately models the behavior of several high-performance disk drives produced in the early 1990s. The same has been found true of disk drives produced in the late 1990s [18].

This manual does not provide details about DiskSim's internals. We refer those that wish to understand DiskSim more thoroughly and/or modify it to the appendices of [2].

#### 1.5.1 Contributors

Many people have contributed to DiskSim's development over the past 11 years including: Bruce Worthington, Steve Schlosser, John Griffin, Ross Cohen, Jiri Schindler, Chris Lumb, John Bucy and Greg Ganger.

## 2 Running DiskSim

DiskSim requires five command line arguments and optionally accepts some number of parameter overrides:

```
disksim <parfile> <outfile> <tracetype> <tracefile> <synthgen> [ par_override ... ]
```

where:

- `disksim` is the name of the executable.
- `parfile` is the name of the parameter file (whose format is described in chapter 3).
- `outfile` is the name of the output file (whose format is described in chapter 5). Output can be directed to stdout by specifying “stdout” for this argument.
- `tracetype` identifies the format of the trace input, if any (options are described in chapter 4).
- `tracefile` identifies the trace file to be used as input. Input is taken from stdin when “stdin” is specified for this argument.
- `synthgen` determines whether or not the synthetic workload generation portion of the simulator should be enabled (any value other than “0” enables synthetic workload generation). The synthetic generator(s) are configured by values in the parameter file, as described in chapter 4. Currently, DiskSim cannot use both an input trace and an internally-generated synthetic workload at the same time.
- `par_override` allows default parameter values or parameter values from `parfile` to be replaced by values specified in the command line. The exact syntax is described in the following section.

### 2.1 Parameter Overrides

When using DiskSim to examine the performance impacts of various storage subsystem design decisions (e.g., sensitivity analyses), the experimental configurations of interest are often quite similar. To avoid the need for numerous parameter files with incremental differences, DiskSim allows parameter values to be overridden with command line arguments. The parameter overrides are applied after the parameter file has been read, but before the internal configuration phase begins. Each parameter override is described by a (component, param name, param value) triple:

```
<component> <parameter> <new value>
```

1. `component` is the name of a component whose parameters are to be overridden. This is the name given to the component when it is instantiated in the parameter file. Ranges are supported; for example `disk0 .. disk5` indicates 6 disks “disk0” through “disk5.” Wildcards are also supported; a trailing `*` matches any string of digits. For example `driver*` matches `driver2`, `driver` and `driver2344` but not `driverqux`.
2. `parameter` is a string identifying the parameter to be overridden. This is identical to the variable name used in the parameter file. If the name contains spaces, it must be quoted so that the shell will pass it to DiskSim as a single argument. To reference a parameter of a subcomponent such as a disk’s scheduler, use the form `Scheduler:parameter`
3. `new value` is the new value for the parameter for the specified instances of the specified module.

Every parameter in the parameter file can be overridden on the command line. Some parameters’ definitions are long enough that it may prove more practical to switch out the parameter files rather than use the command-line override.

## 2.2 Example Command Line

An example may be useful to demonstrate the command line syntax. The following command:

```
disksim parms.1B stdout ascii t.Jan6 0 "disk1 .. disk16" "Segment size (in  
blks)" 64 "disk*" "Scheduler:Scheduling policy" 4
```

executes DiskSim as follows:

- initial parameters are read from file *parms.1B*;
- output (e.g., statistics) is sent to *stdout*;
- the *ascii* input trace is read from file *t.Jan6*;
- there is no synthetically generated activity;
- the cache segment size parameter values of disks 1 through 16, as specified in the parameter file (*parms.1B*), are overridden with a value of 64 (sectors); and
- the scheduling algorithm parameter value for all components matching “disk\*” is overridden with a value of 4 (which corresponds to a Shortest-Seek-Time-First algorithm).



### 3 The Parameter File

DiskSim can be configured via the parameter file to model a wide variety of storage subsystems.

DiskSim uses libparam to input the parameter file; a brief overview of libparam is provided here. At top level in a parameter file, there are three kinds of things in a parameter file: blocks delimited by { }, instantiations, and topology specifications. Instantiations are described in Section 3.5, and topology specifications are described in Section 3.6.

A block consists of a number of “name = value” assignments. Names may contain spaces and are case-sensitive. Values may be integers (including 0x-prefixed hexadecimal), floating point numbers, strings, blocks, and lists delimited by [ ].

Libparam contains a directive called source similar to the #include preprocessor directive in C. source may be used recursively up to a compile-time depth limit of 32.

The organization of the parameter file is quite flexible though there are a few required blocks. Also, components must not be referenced prior to their definition. Every DiskSim parameter file must define the Global and Stats blocks. A simulation using the synthetic trace generator must also define the Proc and Synthio blocks. A typical setup will then define some number of buses, controllers and an iodriver, define or source in some storage device descriptions, instantiate all of these, and then define their interconnection in the simulated storage simulation in a topology specification.<sup>3</sup>

Disk array data organizations are described in logorg blocks (Section 3.8). Every access must fall within some logorg, so at least one must be defined.

Rotational synchronization of devices may optionally be described in a syncset block (Section 3.7).

Adjusting the time scale and remapping requests from a trace may be described in the iosim block (Section 3.3).

Several example parameter files are provided with the software distribution.

The remainder of this section describes each block and its associated parameters. The format of each parameter’s description is:

blockname	paramname	type	required	optional
Parameter description				

which describes a parameter *paramname*, whose type is *type* appearing in block *blockname*.

#### 3.1 Global Block

The global block contains a number of simulation-wide parameters. The name given to the global block to be used must be Global.

disksim_global	Init Seed	int	optional
This specifies the initial seed for the random number generator. The initial seed value is applied at the very beginning of the simulation and is used during the initialization phase (e.g., for determining initial rotational positions). Explicitly specifying the random generator seed enables experiment repeatability.			

disksim_global	Init Seed with time	int	optional
If a nonzero value is provided, DiskSim will use the current system time to initialize the “Init Seed” parameter.			

disksim_global	Real Seed	int	optional
The ‘real’ seed value is applied after the initialization phase and is used during the simulation phase (e.g., for synthetic workload generation). This allows multiple synthetic workloads (with different simulation seeds) to be run on equivalent configurations (i.e., with identical initial seeds, as specified above).			

disksim_global	Real Seed with time	int	optional
If a nonzero value is provided, DiskSim will use the current system time to initialize the “Real Seed” parameter.			

<sup>3</sup>The separation of component definitions and their interconnections greatly reduces the effort required to develop and integrate new components as well as the effort required to understand and modify the existing components [17].

disksim_global	Statistic warm-up time	float	optional
This specifies the amount of simulated time after which the statistics will be reset.			

disksim_global	Statistic warm-up IOs	int	optional
This specifies the number of I/Os after which the statistics will be reset.			

disksim_global	Stat definition file	string	required
This specifies the name of the input file containing the specifications for the statistical distributions to collect. This file allows the user to control the number and sizes of histogram bins into which data are collected. This file is mandatory. Section 5.1 describes its use.			

disksim_global	Output file for trace of I/O requests simulated	string	optional
This specifies the name of the output file to contain a trace of disk request arrivals (in the default ASCII trace format described in Section 4.1). This allows instances of synthetically generated workloads to be saved and analyzed after the simulation completes. This is particularly useful for analyzing (potentially pathological) workloads produced by a system-level model.			

## 3.2 Stats Block

This block contains a series of Boolean [1 or 0] parameters that specify whether or not particular groups of statistics are reported. The name given to the stats block must be `Stats`.

The `iodriver` parameters control statistics output from both the device drivers (individual values and overall totals) and any driver-level disk array logical data organizations (referred to as *logorgs*). The device parameters control statistics output for the devices themselves (individually, overall, and combined with the other devices in a particular logorg). The different print-control parameters (corresponding to particular statistics groups) will be identified with individual statistics in Section 5.

disksim_stats	iodriver stats	block	required
disksim_stats	bus stats	block	required
disksim_stats	ctlr stats	block	required
disksim_stats	device stats	block	required
disksim_stats	process flow stats	block	required

### 3.2.1 Bus Statistics

disksim_bus_stats	Print bus idle stats	int	required
disksim_bus_stats	Print bus arbwait stats	int	required

### 3.2.2 Controller Statistics

disksim_ctlr_stats	Print controller cache stats	int	required
disksim_ctlr_stats	Print controller size stats	int	required
disksim_ctlr_stats	Print controller locality stats	int	required
disksim_ctlr_stats	Print controller blocking stats	int	required
disksim_ctlr_stats	Print controller interference stats	int	required
disksim_ctlr_stats	Print controller queue stats	int	required
disksim_ctlr_stats	Print controller crit stats	int	required
disksim_ctlr_stats	Print controller idle stats	int	required
disksim_ctlr_stats	Print controller intarr stats	int	required

disksim_ctrl_stats	Print controller streak stats	int	required
disksim_ctrl_stats	Print controller stamp stats	int	required
disksim_ctrl_stats	Print controller per-device stats	int	required

### 3.2.3 Device Statistics

disksim_device_stats	Print device queue stats	int	required
disksim_device_stats	Print device crit stats	int	required
disksim_device_stats	Print device idle stats	int	required
disksim_device_stats	Print device intarr stats	int	required
disksim_device_stats	Print device size stats	int	required
disksim_device_stats	Print device seek stats	int	required
disksim_device_stats	Print device latency stats	int	required
disksim_device_stats	Print device xfer stats	int	required
disksim_device_stats	Print device acctime stats	int	required
disksim_device_stats	Print device interfere stats	int	required
disksim_device_stats	Print device buffer stats	int	required

### 3.2.4 iodriver Statistics

disksim_io_driver_stats	Print driver size stats	int	required
disksim_io_driver_stats	Print driver locality stats	int	required
disksim_io_driver_stats	Print driver blocking stats	int	required
disksim_io_driver_stats	Print driver interference stats	int	required
disksim_io_driver_stats	Print driver queue stats	int	required
disksim_io_driver_stats	Print driver crit stats	int	required
disksim_io_driver_stats	Print driver idle stats	int	required
disksim_io_driver_stats	Print driver intarr stats	int	required
disksim_io_driver_stats	Print driver streak stats	int	required
disksim_io_driver_stats	Print driver stamp stats	int	required
disksim_io_driver_stats	Print driver per-device stats	int	required

### 3.2.5 Process-flow Statistics

disksim_pf_stats	Print per-process stats	int	required
disksim_pf_stats	Print per-CPU stats	int	required
disksim_pf_stats	Print all interrupt stats	int	required
disksim_pf_stats	Print sleep stats	int	required

## 3.3 iosim Block

Several aspects of input trace handling are configured in the iosim block.

disksim_iosim	I/O Trace Time Scale	float	optional
<p>This specifies a value by which each arrival time in a trace is multiplied. For example, a value of 2.0 doubles each arrival time, lightening the workload by stretching it out over twice the length of time. Conversely, a value of 0.5 makes the workload twice as heavy by compressing inter-arrival times. This value has no effect on workloads generated internally (by the synthetic generator).</p>			

<code>disksim_iosim</code>	I/O Mappings	list	optional
This is a list of <code>iomap</code> blocks (see below) which enable translation of disk request sizes and locations in an input trace into disk request sizes and locations appropriate for the simulated environment. When the simulated environment closely matches the traced environment, these mappings may be used simply to reassign disk device numbers. However, if the configured environment differs significantly from the trace environment, or if the traced workload needs to be scaled (by request size or range of locations), these mappings can be used to alter the the traced “logical space” and/or scale request sizes and locations. One mapping is allowed per traced device. The mappings from devices identified in the trace to the storage subsystem devices being modeled are provided by block values.			

The I/O Mappings parameter takes a list of `iomap` blocks which contain the following fields:

<code>disksim_iomap</code>	<code>tracedev</code>	int	required
This specifies the traced device affected by this mapping.			

<code>disksim_iomap</code>	<code>simdev</code>	string	required
This specifies the simulated device such requests should access.			

<code>disksim_iomap</code>	<code>locScale</code>	int	required
This specifies a value by which a traced disk request location is multiplied to generate the starting location (in bytes) of the simulated disk request. For example, if the input trace specifies locations in terms of 512-byte sectors, a value of 512 would result in an equivalent logical space of requests.			

<code>disksim_iomap</code>	<code>sizeScale</code>	int	required
This specifies a value by which a traced disk request size is multiplied to generate the size (in bytes) of the simulated disk request.			

<code>disksim_iomap</code>	<code>offset</code>	int	optional
This specifies a value to be added to each simulated request’s starting location. This is especially useful for combining multiple trace devices’ logical space into the space of a single simulated device.			

### 3.4 I/O Subsystem Component Specifications

DiskSim provides four main component types: device drivers, buses, controllers and storage devices. The storage device type currently has three subtypes: a disk model, a simplified, fixed-access-time diskmodel hereinafter referred to as “simpledisk.”

There are two additional components – queues/schedulers and caches – which serve as subcomponents of the above components.

#### 3.4.1 Device Drivers

At the top of the storage system, DiskSim simulates the activity of a device driver in the host system. There must be exactly one configured device driver per storage subsystem configuration.

<code>disksim_iodriver</code>	<code>type</code>	int	required
This is included for extensibility purposes.			

disksim_io_driver	Constant access time	float	required
<p>This specifies any of several forms of storage simulation abstraction. A positive value indicates a fixed access time (after any queueing delays) for each disk request. With this option, requests do not propagate to lower levels of the storage subsystem (and the stats and configuration of lower levels are therefore meaningless). -1.0 indicates that the trace provides a measured access time for each request, which should be used instead of any simulated access times. -2.0 indicates that the trace provides a measured queue time for each request, which should be used instead of any simulated queue times. (Note: This can cause problems if multiple requests are simultaneously issued to disks that don't support queueing.) -3.0 indicates that the trace provides measured values for both the access time and the queue time. Finally, 0.0 indicates that the simulation should compute all access and queue times as appropriate given the changing state of the storage subsystem.</p>			

disksim_io_driver	Use queueing in subsystem	int	required
<p>This specifies whether or not the device driver allows more than one request to be outstanding (in the storage subsystem) at any point in time. During initialization, this parameter is combined with the parameterized capabilities of the subsystem itself to determine whether or not queueing in the subsystem is appropriate.</p>			

disksim_io_driver	Scheduler	block	required
<p>This is an ioqueue; see section 3.4.7.</p>			

### 3.4.2 Buses

disksim_bus	type	int	required
<p>This specifies the type of bus. 1 indicates an exclusively-owned (tenured) bus (i.e., once ownership is acquired, the owner gets 100% of the bandwidth available until ownership is relinquished voluntarily). 2 indicates a shared bus where multiple bulk transfers are interleaved (i.e., each gets a fraction of the total bandwidth).</p>			

disksim_bus	Arbitration type	int	required
<p>This specifies the type of arbitration used for exclusively-owned buses (see above parameter description). 1 indicates slot-based priority (e.g., SCSI buses), wherein the order of attachment determines priority (i.e., the first device attached has the highest priority). 2 indicates First-Come-First-Served (FCFS) arbitration, wherein bus requests are satisfied in arrival order.</p>			

disksim_bus	Arbitration time	float	required
<p>This specifies the time (in milliseconds) required to make an arbitration decision.</p>			

disksim_bus	Read block transfer time	float	required
<p>This specifies the time in milliseconds required to transfer a single 512-byte block in the direction of the device driver / host.</p>			

disksim_bus	Write block transfer time	float	required
<p>This specifies the time (in milliseconds) required to transfer a single 512-byte block in the direction of the disk drives.</p>			

disksim_bus	Print stats	int	required
<p>This specifies whether or not the collected statistics are reported.</p>			

### 3.4.3 Controllers

disksim_ctlr	type	int	required
This specifies the type of controller. 1 indicates a simple controller that acts as nothing more than a bridge between two buses, passing everything straight through to the other side. 2 indicates a very simple, driver-managed controller based roughly on the NCR 53C700. 3 indicates a more complex controller that decouples lower-level storage component peculiarities from higher-level components (e.g., device drivers). The complex controller queues and schedules its outstanding requests and possibly contains a cache. As indicated below, it requires several parameters in addition to those needed by the simpler controllers.			
disksim_ctlr	Scale for delays	float	required
This specifies a multiplicative scaling factor for the various processing delays incurred by the controller. Default overheads for the 53C700-based controller and the more complex controller are hard-coded into the “read_specs” procedure of the controller module (and are easily changed). For the simple pass-thru controller, the scale factor represents the per-message propagation delay (because the hard-coded value is 1.0). 0.0 results in no controller overheads or delays. When the overheads/delays of the controller(s) cannot be separated from those of the disk(s), as is usually the case for single-point tracing of complete systems, the various disk overhead/delay parameter values should be populated and this parameter should be set to 0.0.			
disksim_ctlr	Bulk sector transfer time	float	required
This specifies the time (in milliseconds) necessary to transfer a single 512-byte block to, from or through the controller. Transferring one block over the bus takes the maximum of this time, the block transfer time specified for the bus itself, and the block transfer time specified for the component on the other end of the bus transfer.			
disksim_ctlr	Maximum queue length	int	required
This specifies the maximum number of requests that can be concurrently outstanding at the controller. The device driver discovers this value during initialization and respects it during operation. For the simple types of controllers (see above parameter description), 0 is assumed.			
disksim_ctlr	Print stats	int	required
This specifies whether or not statistics will be reported for the controller. It is meaningless for the simple types of controllers (see above parameter description), as no statistics are collected.			
disksim_ctlr	Scheduler	block	optional
This is an iqueue; see section 3.4.7			
disksim_ctlr	Cache	block	optional
A block cache; see section 3.4.8			
disksim_ctlr	Max per-disk pending count	int	optional
This specifies the maximum number of requests that the controller can have outstanding to each attached disk (i.e., the maximum number of requests that can be dispatched to a single disk). This parameter only affects the interaction of the controller with its attachments; it is not visible to the device driver.			

### 3.4.4 Storage Devices

“Storage devices” are the abstraction through which the various storage device models are interfaced with disksim. In the current release, there are 2 such models: conventional disks, and a simplified, fixed-access-time diskmodel hereinafter referred to as “simpldisk.”

### 3.4.5 Disks

Since disk specifications are long, it is often convenient to store them in separate files and include them in the main parameter file via the “source” directive.

disksim_disk	Model	block	required
Parameters for the disk's libdiskmodel model. See the diskmodel documentation for details.			
disksim_disk	Scheduler	block	required
An ioqueue; see Section 3.4.7			
disksim_disk	Max queue length	int	required
This specifies the maximum number of requests that the disk can have in service or queued for service at any point in time. During initialization, other components request this information and respect it during simulation.			
disksim_disk	Bulk sector transfer time	float	required
This specifies the time for the disk to transfer a single 512-byte block over the bus. Recall that this value is compared with the corresponding bus and controller block transfer values to determine the actual transfer time (i.e., the maximum of the three values).			
disksim_disk	Segment size (in blks)	int	required
This specifies the size of each buffer segment, assuming a static segment size. Some modern disks will dynamically resize their buffer segments (and thereby alter the number of segments) to respond to perceived patterns of workload behavior, but DiskSim does not currently support this functionality.			
disksim_disk	Number of buffer segments	int	required
This specifies the number of segments in the on-board buffer/cache. A buffer segment is similar to a cache line, in that each segment contains data that is disjoint from all other segments. However, segments tend to be organized as circular queues of logically sequential disk sectors, with new sectors pushed into an appropriate queue either from the bus (during a write) or from the disk media (during a read). As data are read from the buffer/cache and either transferred over the bus (during a read) or written to the disk media (during a write), they are eligible to be pushed out of the segment (if necessary or according to the dictates of the buffer/cache management algorithm).			
disksim_disk	Print stats	int	required
This specifies whether or not statistics for the disk will be reported.			
disksim_disk	Per-request overhead time	float	required
This specifies a per-request processing overhead that takes place immediately after the arrival of a new request at the disk. It is additive with various other processing overheads described below, but in general either the other overheads are set to zero or this parameter is set to zero.			
disksim_disk	Time scale for overheads	float	required
This specifies a multiplicative scaling factor for all processing overhead times. For example, 0.0 eliminates all such delays, 1.0 uses them at face value, and 1.5 increases them all by 50%.			
disksim_disk	Hold bus entire read xfer	int	required
This specifies whether or not the disk retains ownership of the bus throughout the entire transfer of “read” data from the disk. If false (0), the disk may release the bus if and when the current transfer has exhausted all of the available data in the on-board buffer/cache and must wait for additional data sectors to be read off the disk into the buffer/cache.			

disksim_disk	Hold bus entire write xfer	int	required
This specifies whether or not the disk retains ownership of the bus throughout the entire transfer of “write” data to the disk. If false (0), the disk may release the bus if and when the current transfer has filled up the available space in the on-board buffer/cache and must wait for data from the buffer/cache to be written to the disk.			
disksim_disk	Allow almost read hits	int	required
This specifies whether or not a new read request whose first block is currently being prefetched should be treated as a partial cache hit. Doing so generally means that the request is handled right away.			
disksim_disk	Allow sneaky full read hits	int	required
This specifies whether or not a new read request whose data are entirely contained in a single segment of the disk cache is allowed to immediately transfer that data over the bus while another request is moving the disk actuator and/or transferring data between the disk cache and the disk media. In essence, the new read request “sneaks” its data out from the disk cache without interrupting the current (active) disk request.			
disksim_disk	Allow sneaky partial read hits	int	required
This specifies whether or not a new read request whose data are partially contained in a single segment of the disk cache is allowed to immediately transfer that data over the bus while another request is moving the disk actuator and/or transferring data between the disk cache and the disk media. In essence, the new read request “sneaks” the cached portion of its data out from the disk cache without interrupting the current (active) disk request.			
disksim_disk	Allow sneaky intermediate read hits	int	required
This specifies whether or not the on-board queue of requests is searched during idle bus periods in order to find read requests that may be partially or completely serviced from the current contents of the disk cache. That is, if the current (active) request does not need bus access at the current time, and the bus is available for use, a queued read request whose data are in the cache may obtain access to the bus and begin data transfer. “Full” intermediate read hits are given precedence over “partial” intermediate read hits.			
disksim_disk	Allow read hits on write data	int	required
This specifies whether or not data placed in the disk cache by write requests are considered usable by read requests. If false (0), such data are removed from the cache as soon as they have been copied to the media.			
disksim_disk	Allow write prebuffering	int	required
This specifies whether or not the on-board queue of requests is searched during idle bus periods for write requests that could have part or all of their data transferred to the on-board cache (without disturbing an ongoing request). That is, if the current (active) request does not need bus access at the current time, and the bus is available for use, a queued write request may obtain access to the bus and begin data transfer into an appropriate cache segment. Writes that are contiguous to the end of the current (active) request are given highest priority in order to facilitate continuous transfer to the media, followed by writes that have already “prebuffered” some portion of their data.			
disksim_disk	Preseeking level	int	required
This specifies how soon the actuator is allowed to start seeking towards the media location of the next request’s data. 0 indicates no preseeking, meaning that the actuator does not begin relocation until the current request’s completion has been confirmed by the controller (via a completion “handshake” over the bus). 1 indicates that the actuator can begin relocation as soon as the completion message has been prepared for transmission by the disk. 2 indicates that the actuator can begin relocation as soon as the access of the last sector of the current request (and any required prefetching) has been completed. This allows greater parallelism between bus activity and mechanical activity.			



disksim_disk	Never disconnect	int	required
This specifies whether or not the disk retains ownership of the bus during the entire processing and servicing of a request (i.e., from arrival to completion). If false (0), the disk may release the bus whenever it is not needed for transferring data or control information.			
disksim_disk	Avg sectors per cylinder	int	required
This specifies (an approximation of) the mean number of sectors per cylinder. This value is exported to external schedulers some of which utilize approximations of the actual layout of data on the disk(s) when reordering disk requests. This value is not used by the disk itself.			
disksim_disk	Maximum number of write segments	int	required
This specifies the number of cache segments available for holding “write” data at any point in time. Because write-back caching is typically quite limited in current disk cache management schemes, some caches only allow a subset of the segments to be used to hold data for write requests (in order to minimize any interference with sequential read streams).			
disksim_disk	Use separate write segment	int	required
This specifies whether or not a single segment is statically designated for use by all write requests. This further minimizes the impact of write requests on the caching/prefetching of sequential read streams.			
disksim_disk	Low (write) water mark	float	required
This specifies the fraction of segment size or request size (see below) corresponding to the <i>low water mark</i> . When data for a write request are being transferred over the bus into the buffer/cache, and the buffer/cache segment fills up with “dirty” data, the disk may disconnect from the bus while the buffered data are written to the disk media. When the amount of dirty data in the cache falls below the low water mark, the disk attempts to reconnect to the bus to continue the interrupted data transfer.			
disksim_disk	High (read) water mark	float	required
This specifies the fraction of segment size or request size (see below) corresponding to the <i>high water mark</i> . When data for a read request are being transferred over the bus from the buffer/cache, and the buffer/cache segment runs out of data to transfer, the disk may disconnect from the bus until additional data are read from the disk media. When the amount of available data in the cache reaches the high water mark, the disk attempts to reconnect to the bus to continue the interrupted data transfer.			
disksim_disk	Set watermark by reqsize	int	required
This specifies whether the watermarks are computed as fractions of the individual request size or as fractions of the buffer/cache segment size.			
disksim_disk	Calc sector by sector	int	required
This specifies whether or not media transfers should be computed sector by sector rather than in groups of sectors. This optimization has no effect on simulation accuracy, but potentially results in shorter simulation times (at a cost of increased code complexity). It has not been re-enabled since the most recent modifications to DiskSim, so the simulator currently functions as if the value were always true (1).			
disksim_disk	Enable caching in buffer	int	required
This specifies whether or not on-board buffer segments are used for data caching as well as for speed-matching between the bus and the disk media. Most (if not all) modern disk drives utilize their buffers as caches.			

disksim_disk	Buffer continuous read	int	required
<p>This specifies the type of prefetching performed by the disk. 0 disables prefetching. 1 enables prefetching up to the end of the track containing the last sector of the read request. 2 enables prefetching up to the end of the cylinder containing the last sector of the read request. 3 enables prefetching up to the point that the current cache segment is full. 4 enables prefetching up to the end of the track following the track containing the last sector of the read request, provided that the current request was preceded in the not-too-distant past by another read request that accessed the immediately previous track. In essence, the last scheme enables a type of prefetching that tries to stay one logical track “ahead” of any sequential read streams that are detected.</p>			
disksim_disk	Minimum read-ahead (blks)	int	required
<p>This specifies the minimum number of disk sectors that must be prefetched after a read request before servicing another (read or write) request. A positive value may be beneficial for workloads containing multiple interleaved sequential read streams, but 0 is typically the appropriate value.</p>			
disksim_disk	Maximum read-ahead (blks)	int	required
<p>This specifies the maximum number of disk sectors that may be prefetched after a read request (regardless of all other prefetch parameters).</p>			
disksim_disk	Read-ahead over requested	int	required
<p>This specifies whether or not newly prefetched data can replace (in a buffer segment) data returned to the host as part of the most recent read request.</p>			
disksim_disk	Read-ahead on idle hit	int	required
<p>This specifies whether or not prefetching should be initiated by the disk when a read request is completely satisfied by cached data (i.e., a “full read hit”).</p>			
disksim_disk	Read any free blocks	int	required
<p>This specifies whether or not disk sectors logically prior to the requested sectors should be read into the cache if they pass under the read/write head prior to reaching the requested data (e.g., during rotational latency).</p>			
disksim_disk	Fast write level	int	required
<p>This specifies the type of write-back caching implemented. 0 indicates that write-back caching is disabled (i.e., all dirty data must be written to the disk media prior to sending a completion message). 1 indicates that write-back caching is enabled for contiguous sequential write request streams. That is, as long as each request arriving at the disk is a write request that “appends” to the current segment of dirty data, a completion message will be returned for each new request as soon as all of its data have been transferred over the bus to the disk buffer/cache. 2 indicates that write-back caching is enabled for contiguous sequential write request streams even if they are intermixed with read or non-appending write requests, although before any such request is serviced by the disk, all of the dirty write data must be flushed to the media. A scheduling algorithm that gives precedence to sequential writes would maximize the effectiveness of this option.</p>			
disksim_disk	Combine seq writes	int	required
<p>This specifies whether or not sequential data from separate write requests can share a common cache segment. If true (1), data are typically appended at the end of a previous request’s dirty data. However, if all of the data in a cache segment are dirty, and no mechanical activity has begun on behalf of the request(s) using that segment, “prepending” of additional dirty data are allowed provided that the resulting cache segment contains a single contiguous set of dirty sectors.</p>			
disksim_disk	Stop prefetch in sector	int	required
<p>This specifies whether or not a prefetch may be aborted in the “middle” of reading a sector off the media. If false (0), prefetch activity is only aborted at sector boundaries.</p>			

disksim_disk	Disconnect write if seek	int	required
This specifies whether or not the disk should disconnect from the bus if the actuator is still in motion (seeking) when the last of a write request's data has been transferred to the disk buffer/cache.			
disksim_disk	Write hit stop prefetch	int	required
This specifies whether or not the disk should discontinue the read-ahead of a previous request when a write hit in the cache occurs. Doing so allows the new write request's data to begin travelling to the disk more quickly, at the expense of some prefetching activity.			
disksim_disk	Read directly to buffer	int	required
This specifies whether or not space for a sector must be available in the buffer/cache prior to the start of the sector read. If false (0), a separate sector buffer is assumed to be available for use by the media-reading electronics, implying that the data for a sector is transferred to the main buffer/cache only after it has been completely read (and any error-correction algorithms completed).			
disksim_disk	Immed transfer partial hit	int	required
This specifies whether or not a read request whose initial (but not all) data are present in the disk buffer/cache has that data immediately transferred over the bus. If false (0), the data are immediately transferred only if the amount of requested data that are present in the buffer/cache exceed the <i>high water mark</i> (see above).			
disksim_disk	Read hit over. after read	float	required
This specifies the processing time for a read request that hits in the on-board cache when the immediately previous request was also a read. This delay is applied before any data are transferred from the disk buffer/cache.			
disksim_disk	Read hit over. after write	float	required
This specifies the processing time for a read request that hits in the on-board cache when the immediately previous request was a write. This delay is applied before any data are transferred from the disk buffer/cache.			
disksim_disk	Read miss over. after read	float	required
This specifies the processing time for a read request that misses in the on-board cache when the immediately previous request was also a read. This delay is applied before any mechanical positioning delays or data transfer from the media.			
disksim_disk	Read miss over. after write	float	required
This specifies the processing time for a read request that misses in the on-board cache when the immediately previous request was a write. This delay is applied before any mechanical positioning delays or data transfer from the media.			
disksim_disk	Write hit over. after read	float	required
This specifies the processing time for a write request that "hits" in the on-board cache (i.e., completion will be reported before data reaches media) when the immediately previous request was a read. This delay is applied before any mechanical positioning delays and before any data are transferred to the disk buffer/cache.			
disksim_disk	Write hit over. after write	float	required
This specifies the processing time for a write request that "hits" in the on-board cache (i.e., completion will be reported before data reaches media) when the immediately previous request was also a write. This delay is applied before any mechanical positioning delays and before any data are transferred to the disk buffer/cache.			

disksim_disk	Write miss over. after read	float	required
This specifies that “misses” in the on-board cache (i.e., completion will be reported only after data reaches media) when the immediately previous request was a read. This delay is applied before any mechanical positioning delays and before any data are transferred to the disk buffer/cache.			
disksim_disk	Write miss over. after write	float	required
This specifies the processing time for a write request that “misses” in the on-board cache (i.e., completion will be reported only after data reaches media) when the immediately previous request was also a write. This delay is applied before any mechanical positioning delays and before any data are transferred to the disk buffer/cache.			
disksim_disk	Read completion overhead	float	required
This specifies the processing time for completing a read request. This overhead is applied just before the completion message is sent over the (previously acquired) bus and occurs in parallel with any background disk activity (e.g., prefetching or preseeking).			
disksim_disk	Write completion overhead	float	required
This specifies the processing time for completing a write request. This overhead is applied just before the completion message is sent over the (previously acquired) bus and occurs in parallel with any background disk activity (e.g., preseeking).			
disksim_disk	Data preparation overhead	float	required
This specifies the additional processing time necessary when preparing to transfer data over the bus (for either reads or writes). This command processing overhead is applied after obtaining access to the bus (prior to transferring any data) and occurs in parallel with any ongoing media access.			
disksim_disk	First reselect overhead	float	required
This specifies the processing time for reconnecting to (or “reselecting”) the controller for the first time during the current request. This command processing overhead is applied after the disk determines that reselection is appropriate (prior to attempting to acquire the bus) and occurs in parallel with any ongoing media access. Reselection implies that the disk has explicitly disconnected from the bus at some previous point while servicing the current request and is now attempting to reestablish communication with the controller. Disconnection and subsequent reselection result in some additional command processing and protocol overhead, but they may also improve the overall utilization of bus resources shared by multiple disks (or other peripherals).			
disksim_disk	Other reselect overhead	float	required
This specifies the processing time for reconnecting to the controller after the first time during the current request (i.e., the second reselection, the third reselection, etc.). This command processing overhead is applied after the disk determines that reselection is appropriate (prior to attempting to acquire the bus) and occurs in parallel with any ongoing media access.			
disksim_disk	Read disconnect afterread	float	required
This specifies the processing time for a read request that disconnects from the bus when the previous request was also a read. This command processing overhead is applied after the disk determines that disconnection is appropriate (prior to requesting disconnection from the bus) and occurs in parallel with any ongoing media access.			

disksim_disk	Read disconnect afterwrite	float	required
This specifies the processing time for a read request that disconnects from the bus when the previous request was a write request. This command processing overhead is applied after the disk determines that disconnection is appropriate (prior to requesting disconnection from the bus) and occurs in parallel with any ongoing media access.			
disksim_disk	Write disconnect overhead	float	required
This specifies the processing time for a write request that disconnects from the bus (which generally occurs after the data are transferred from the host to the on-board buffer/cache). This command processing overhead is applied after the disk determines that disconnection is appropriate (prior to requesting disconnection from the bus) and occurs in parallel with any ongoing media access.			
disksim_disk	Extra write disconnect	int	required
This specifies whether or not the disk disconnects from the bus after processing the write command but before any data have been transferred over the bus into the disk buffer/cache. Although there are no performance or reliability advantages to this behavior, it has been observed in at least one production SCSI disk and has therefore been included in DiskSim. If true (1), the next five parameters configure additional overhead values specifically related to this behavior.			
disksim_disk	Extradisc command overhead	float	required
This specifies the processing time for a write request that disconnects from the bus before transferring any data to the disk buffer/cache. This overhead is applied before requesting disconnection from the bus and before any mechanical positioning delays. This parameter (when enabled) functions in place of the above "Write over." parameters.			
disksim_disk	Extradisc disconnect overhead	float	required
This specifies the additional processing time for a write request that disconnects from the bus before transferring any data to the disk buffer/cache. This overhead is also applied before requesting disconnection from the bus, but it occurs in parallel with any mechanical positioning delays. This parameter (when enabled) functions in place of the above "Write disconnect" parameter for initial write disconnections.			
disksim_disk	Extradisc inter-disconnect delay	float	required
This specifies the time between the initial disconnect from the bus and the subsequent reconnection attempt for a write request that disconnects from the bus before transferring any data to the disk buffer/cache. It occurs in parallel with any mechanical positioning delays.			
disksim_disk	Extradisc 2nd disconnect overhead	float	required
This specifies the processing time for a write request that disconnects from the bus after data has been transferred but previously had disconnected without transferring any data to the disk buffer/cache. This command processing overhead is applied after the disk determines that disconnection is appropriate (prior to requesting disconnection from the bus) and occurs in parallel with any ongoing media access. This parameter (when enabled) functions in place of the above "Write disconnect" parameter for non-initial write disconnections.			
disksim_disk	Extradisc seek delta	float	required
This specifies the additional delay between the completion of the initial command processing overhead and the initiation of any mechanical positioning for a write request that disconnects from the bus before transferring any data to the disk buffer/cache. This delay occurs in parallel with ongoing bus activity and related processing overheads.			

disksim_disk	Minimum seek delay	float	required
This specifies the minimum media access delay for a nonsequential write request. That is, a nonsequential write request (after any command processing overheads) must wait at least this amount of time before accessing the disk media.			

disksim_disk	Immediate buffer read	int	required
This specifies whether or not disk sectors should be transferred into the on-board buffer in the order that they pass under the read/write head rather than in strictly ascending logical block order. This is known as <i>zero-latency reads</i> or <i>read-on-arrival</i> . It is intended to improve response times by reducing rotational latency (by not rotating all the way around to the first requested sector before beginning to fill the buffer/cache).			

disksim_disk	Immediate buffer write	int	required
This specifies whether or not disk sectors should be transferred from the on-board buffer in the order that they pass under the read/write head rather than in strictly ascending logical block order. These are known as <i>zero-latency writes</i> or <i>write-on-arrival</i> . It is intended to improve response times by reducing rotational latency (by not rotating all the way around to the first “dirty” sector before beginning to flush the buffer/cache).			

### 3.4.6 Simple Disks

The `simpledisk` module provides a simplified model of a storage device that has a constant access time. It was implemented mainly as an example and to test the interface through which new storage device types might later be added to `DiskSim`.

disksim_simpledisk	Scheduler	block	required
This is an ioqueue; see Section 3.4.7 for details.			

disksim_simpledisk	Max queue length	int	required
This specifies the maximum number of requests that the <code>simpledisk</code> can have in service or queued for service at any point in time. During initialization, other components request this information and respect it during simulation.			

disksim_simpledisk	Block count	int	required
This specifies the capacity of the <code>simpledisk</code> in blocks.			

disksim_simpledisk	Bus transaction latency	float	optional
This specifies the delay involved at the <code>simpledisk</code> for each message that it transfers.			

disksim_simpledisk	Bulk sector transfer time	float	required
This specifies the time necessary to transfer a single 512-byte block to, from or through the controller. Transferring one block over the bus takes the maximum of this time, the block transfer time specified for the bus itself, and the block transfer time specified for the component on the other end of the bus transfer.			

disksim_simpledisk	Never disconnect	int	required
This specifies whether or not the <code>simpledisk</code> retains ownership of the bus during the entire processing and servicing of a request (i.e., from arrival to completion). If false (0), the <code>simpledisk</code> may release the bus whenever it is not needed for transferring data or control information.			

disksim_simpledisk	Print stats	int	required
Specifies whether or not statistics for the <code>simpledisk</code> will be reported.			

disksim_simpledisk	Command overhead	float	required
This specifies a per-request processing overhead that takes place immediately after the arrival of a new request at the disk.			
disksim_simpledisk	Constant access time	float	optional
This specifies the fixed per-request access time (i.e., actual mechanical activity is not simulated).			
disksim_simpledisk	Access time	float	required
Synonym for Constant access time.			

### 3.4.7 Queue/Scheduler Subcomponents

disksim_ioqueue	Scheduling policy	int	required
This specifies the primary scheduling algorithm employed for selecting the next request to be serviced. A large set of algorithms have been implemented, ranging from common choices like First-Come-First-Served (FCFS) and Shortest-Seek-Time-First (SSTF) to new algorithms like Shortest-Positioning-(w/Cache)-Time-First (described in [22]). See Table 1 for the list of algorithms provided.			
disksim_ioqueue	Cylinder mapping strategy	int	required
This specifies the level of detail of physical data layout information available to the scheduler. 0 indicates that the only information available to the scheduler are the logical block numbers specified in the individual requests. 1 indicates that the scheduler has access to information about zone boundaries, the number of physical sectors/zone, and the number of physical sectors/track in each zone. 2 indicates that the scheduler also has access to the layout of spare sectors or tracks in each zone. 3 indicates that the scheduler also has access to the list of any slipped sectors/tracks. 4 indicates that the scheduler also has access to the list of any remapped sectors/tracks, thereby providing an exact data layout (logical-to-physical mapping) for the disk. 5 indicates that the scheduler uses the cylinder number given to it with the request, allowing experiments with arbitrary mappings. In particular, some traces include the cylinder number as part of the request record. 6 indicates that the scheduler only has access to (an approximation of) the mean number of sectors per cylinder. The value used in this case is that specified in the disk parameter "Avg. sectors per cylinder."			
disksim_ioqueue	Write initiation delay	float	required
This specifies an approximation of the write request processing overhead(s) performed prior to any mechanical positioning delays. This value is used by scheduling algorithms that select the order of request service based (at least in part) on expected positioning delays.			
disksim_ioqueue	Read initiation delay	float	required
This specifies an approximation of the read request processing overhead performed prior to any mechanical positioning delays. This value is used by scheduling algorithms that select the order of request service based (at least in part) on expected positioning delays.			
disksim_ioqueue	Sequential stream scheme	int	required
The integer value is interpreted as a boolean bitfield. It specifies the type of sequential stream detection and/or request concatenation performed by the scheduler (see [21] for additional details). Bit 0 indicates whether or not sequential read requests are concatenated by the scheduler. Bit 1 indicates whether or not sequential write requests are concatenated by the scheduler. Bit 2 indicates whether or not sequential read requests are always scheduled together. Bit 3 indicates whether or not sequential write requests are always scheduled together. Bit 4 indicates whether or not sequential requests of any kind (e.g., interleaved reads and writes) are always scheduled together.			

disksim_ioqueue	Maximum concat size	int	required
This specifies the maximum request size resulting from concatenation of sequential requests. That is, if the sum of the sizes of the two requests to be concatenated exceeds this value, the concatenation will not be performed by the scheduler.			

disksim_ioqueue	Overlapping request scheme	int	required
This specifies the scheduler's policy for dealing with overlapping requests. 0 indicates that overlapping requests are treated as independent. 1 indicates that requests that are completely overlapped by a completed request that arrived after them are subsumed by that request. 2 augments this policy by also allowing read requests that arrive after the completed overlapping request to be subsumed by it, since the necessary data are known. This support was included for experiments in [2] in order to partially demonstrate the correctness problems of open workloads (e.g., feedback-free request traces). In most real systems, overlapping requests are almost never concurrently outstanding.			

disksim_ioqueue	Sequential stream diff maximum	int	required
This specifies the maximum distance between two "sequential" requests in a sequential stream. This allows requests with a small stride or an occasional "skip" to still be considered for inclusion in a sequential stream.			

disksim_ioqueue	Scheduling timeout scheme	int	required
This specifies the type of multi-queue timeout scheme implemented. 0 indicates that requests are not moved from the <i>base</i> queue to a higher-priority queue because of excessive queueing delays. 1 indicates that requests in the base queue whose queueing delays exceed the specified timeout value (see below) will be moved to one of two higher-priority queues (the <i>timeout</i> queue or the <i>priority</i> queue) based on the scheduling priority scheme (see below). 2 indicates that requests in the base queue whose queueing delays exceed half of the specified timeout value (see below) will be moved to the next higher priority queue (the <i>timeout</i> queue). Furthermore, such requests will be moved to the highest priority queue (the <i>priority</i> queue) if their total queueing delays exceed the specified timeout value (see below).			

disksim_ioqueue	Timeout time/weight	int	required
This specifies either the timeout value (in seconds) for excessive queueing delays or the time/aging factor used in calculating request priorities for various age-sensitive scheduling algorithms. The time/aging factor is additive for some algorithms and multiplicative for others.			

disksim_ioqueue	Timeout scheduling	int	required
This specifies the scheduling algorithm employed for selecting the next request to be serviced from the <i>timeout</i> queue. The options are the same as those available for the "Scheduling policy" parameter above.			

disksim_ioqueue	Scheduling priority scheme	int	required
This specifies whether or not requests flagged as high priority (i.e., time-critical or time-limited requests [5]) are automatically placed in the highest-priority queue (the <i>priority</i> queue).			

disksim_ioqueue	Priority scheduling	int	required
This specifies the scheduling algorithm employed for selecting the next request to be serviced from the <i>priority</i> queue. The options are the same as those available for the "Scheduling policy" parameter above.			

### 3.4.8 Disk Block Cache Subcomponents

The following parameters configure the generic disk block cache subcomponent, which is currently only used in the intelligent controller type (3). The disk module has its own cache submodule, which is configured by disk parameters described in Section 3.4.5.



FCFS	1
ELEVATOR_LBN	2
CYCLE_LBN	3
SSTF_LBN	4
ELEVATOR_CYL	5
CYCLE_CYL	6
SSTF_CYL	7
SPTF_OPT	8
SPCTF_OPT	9
SATF_OPT	10
WPTF_OPT	11
WPCTF_OPT	12
WATF_OPT	13
ASPTF_OPT	14
ASPCTF_OPT	15
ASATF_OPT	16
VSCAN_LBN	17
VSCAN_CYL	18
PRI_VSCAN_LBN	19
PRI_ASPTF_OPT	20
PRI_ASPCTF_OPT	21
SDF_APPROX	22
SDF_EXACT	23
SPTF_ROT_OPT	24
SPTF_ROT_WEIGHT	25
SPTF_SEEK_WEIGHT	26
TSPS	27

Table 1: Scheduling algorithms provided by ioqueues

### 3.4.9 Memory Caches

disksim_cachemem	Cache size	int	required
This specifies the total size of the cache in blocks.			
disksim_cachemem	SLRU segments	list	optional
This is a list of segment sizes in blocks for the segmented-LRU algorithm [10] if it is specified as the cache replacement algorithm (see below).			
disksim_cachemem	Line size	int	required
This specifies the cache line size (i.e., the unit of cache space allocation/replacement).			
disksim_cachemem	Bit granularity	int	required
This specifies the number of blocks covered by each “present” bit and each “dirty” bit. The value must divide the cache line size evenly. Higher values (i.e., coarser granularities) can result in increased numbers of installation reads (i.e., fill requests required to complete partial-line writes [13]).			
disksim_cachemem	Lock granularity	int	required
This specifies the number of blocks covered by each lock. The value must divide the cache line size evenly. Higher values (i.e., coarser granularities) can lead to increased lock contention.			
disksim_cachemem	Shared read locks	int	required
This specifies whether or not read locks are sharable. If false (0), read locks are exclusive.			
disksim_cachemem	Max request size	int	required
This specifies the maximum request size to be served by the cache. This value does not actually affect the simulated cache’s behavior. Rather, higher-level system components (e.g., the device driver in DiskSim) acquire this information at initialization time and break up larger requests to accommodate it. 0 indicates that there is no maximum request size.			
disksim_cachemem	Replacement policy	int	required
This specifies the line replacement policy. 1 indicates First-In-First-Out (FIFO). 2 indicates segmented-LRU [10]. 3 indicates random replacement. 4 indicates Last-In-First-Out (LIFO).			
disksim_cachemem	Allocation policy	int	required
This specifies the line allocation policy. 0 indicates that the cache replacement policy is strictly followed; if the selected line is dirty, the allocation waits for the required write-back request to complete. 1 indicates that “clean” lines are considered for replacement prior to “dirty” lines (and background write-back requests are issued for each dirty line considered).			
disksim_cachemem	Write scheme	int	required
This specifies the policy for handling write requests. 1 indicates that new data are always synchronously written to the backing store before indicating completion. 2 indicates a write-through scheme where requests are immediately initiated for writing out the new data to the backing store. The original write requests are considered complete as soon as the new data is cached. 3 indicates a write-back scheme where completions are reported immediately and dirty blocks are held in the cache for some time before being written out to the backing store.			

disksim_cachemem	Flush policy	int	required
This specifies the policy for flushing dirty blocks to the backing store (assuming a write-back scheme for handling write requests). 0 indicates that dirty blocks are written back “on demand” (i.e., only when the allocation/replacement policy needs to reclaim them). 1 indicates write-back requests are periodically initiated for all dirty cache blocks.			
disksim_cachemem	Flush period	float	required
This specifies the time between periodic write-backs of all dirty cache blocks (assuming a periodic flush policy).			
disksim_cachemem	Flush idle delay	float	required
This specifies the amount of contiguous idle time that must be observed before background write-backs of dirty cache blocks are initiated. Any front-end request processing visible to the cache resets the idle timer. -1.0 indicates that idle background flushing is disabled.			
disksim_cachemem	Flush max line cluster	int	required
This specifies the maximum number of cache lines that can be combined into a single write-back request (assuming “gather” write support).			
disksim_cachemem	Read prefetch type	int	required
This specifies the prefetch policy for handling read requests. Prefetching is currently limited to extending requested fill accesses to include other portions of requested lines. 0 indicates that prefetching is disabled. 1 indicates that unrequested data at the start of a requested line are prefetched. 2 indicates that unrequested data at the end of a requested line are prefetched. 3 indicates that any unrequested data in a requested line are prefetched (i.e., full line fills only).			
disksim_cachemem	Write prefetch type	int	required
This specifies the prefetch policy for handling installation reads (caused by write requests). Prefetching is currently limited to extending the requested fill accesses to include other portions of the requested lines. 0 indicates that prefetching is disabled. 1 indicates that unrequested data at the start of a requested line are prefetched. 2 indicates that unrequested data at the end of a requested line are prefetched. 3 indicates that any unrequested data in a requested line are prefetched (i.e., full line fills only).			
disksim_cachemem	Line-by-line fetches	int	required
This specifies whether or not every requested cache line results in a separate fill request. If false (0), multi-line fill requests can be generated when appropriate.			
disksim_cachemem	Max gather	int	required
This specifies the maximum number of non-contiguous cache lines (in terms of their memory addresses) that can be combined into a single disk request, assuming that they correspond to contiguous disk addresses. (DiskSim currently treats every pair of cache lines as non-contiguous in memory.) 0 indicates that any number of lines can be combined into a single request (i.e., there is no maximum).			

### 3.4.10 Cache Devices

DiskSim can use one device as a cache for another.

disksim_cachedev	Cache size	int	required
This specifies the total size of the cache in blocks.			

disksim_cachedev	Max request size	int	required
This specifies the maximum request size to be served by the cache. This value does not actually affect the simulated cache's behavior. Rather, higher-level system components (e.g., the device driver in DiskSim) acquire this information at initialization time and break up larger requests to accommodate it. 0 indicates that there is no maximum request size.			
disksim_cachedev	Write scheme	int	required
This specifies the policy for handling write requests. 1 indicates that new data are always synchronously written to the backing store before indicating completion. 2 indicates a write-through scheme where requests are immediately initiated for writing out the new data to the backing store, but the original write requests are considered complete as soon as the new data is cached. 3 indicates a write-back scheme where completions are reported immediately and dirty blocks are held in the cache for some time before being written out to the backing store.			
disksim_cachedev	Flush policy	int	required
This specifies the policy for flushing dirty blocks to the backing store (assuming a write-back scheme for handling write requests). 0 indicates that dirty blocks are written back "on demand" (i.e., only when the allocation/replacement policy needs to reclaim them). 1 indicates write-back requests are periodically initiated for all dirty cache blocks.			
disksim_cachedev	Flush period	float	required
This specifies the time between periodic write-backs of all dirty cache blocks (assuming a periodic flush policy).			
disksim_cachedev	Flush idle delay	float	required
This specifies the amount of contiguous idle time that must be observed before background write-backs of dirty cache blocks are initiated. Any front-end request processing visible to the cache resets the idle timer. -1.0 indicates that idle background flushing is disabled.			
disksim_cachedev	Cache device	string	required
The device used for the cache.			
disksim_cachedev	Cached device	string	required
The device whose data is being cached.			

### 3.5 Component Instantiation

The input component specifications must be instantiated and given names before they can be incorporated into a simulated storage system. Component instantiations have the following form:

```
instantiate <name list> as <instance name>
```

where <instance name> is the name given to the component specification and <name list> is a list of names for the instantiated devices.

e.g. `instantiate [ bus0 ] as BUS0`  
creates a bus named `bus0` using the `BUS0` specification.

`instantiate [ disk0, disk2, disk4 .. disk6 ] as IBM.DNES-309170W.validate`  
creates 5 disks with names `disk0`, `disk2`, `disk4`, `disk5` and `disk6` using the `IBM.DNES-309170W.validate` specification.

### 3.6 I/O Subsystem Interconnection Specifications

The allowed interconnections are independent of the components themselves except that a device driver must be at the “top” of any subsystem and storage devices must be at the “bottom.” Exactly one or two controllers must be between the device driver and each disk, with a bus connecting each such pair of components along the path from driver to disk. Each disk or controller can only be connected to one bus from the host side of the subsystem. A bus can have no more than 15 disks or controllers attached to it. A controller can have no more than 4 back-end buses (use of more than one is not well tested). The one allowable device driver is connected to the top-level bus.

The system topology is specified to DiskSim via a topology specification. A topology specification consists of a device type, a device name and a list of devices which are children of that device. The named devices must be instantiated; the component instantiations should precede the topology specification in the parameter file. In the current implementation, no device may appear more than once in the topology specification. Future versions may provide multi-path support.

An example topology specification is provided in Figure 1 below along with a diagram of the storage system corresponding to the specification (Figure 2).

### 3.7 Rotational Synchronization of Devices

DiskSim can be configured to simulate rotationally synchronized devices via the following parameters. Rotationally synchronized devices are always at exactly the same rotational offset, which requires that they begin the simulation at the same offset and rotate at the same speed. Non-synchronized devices are assigned a random initial rotational offset at the beginning of the simulation and are individually assigned a rotational speed based on the appropriate device parameters.

disksim_syncset	type	string	required
The type of devices appearing in the syncset. Currently, only “disk” is allowed.			

disksim_syncset	devices	list	required
A list of names of devices that are in the syncset.			

### 3.8 Disk Array Data Organizations

DiskSim can simulate a variety of logical data organizations, including striping and various RAID architectures. Although DiskSim is organized to allow such organizations both at the system-level (i.e., at front end of the device drivers) and at the controller-level, only system-level organizations are supported in the first released version. Each logical organization is configured in a “logorg” block.

disksim_logorg	Addressing mode	string	required
This specifies how the logical data organization is addressed. <code>Array</code> indicates that there is a single logical device number for the entire logical organization. <code>Parts</code> indicates that back-end storage devices are addressed as though there were no logical organization, and requests are re-mapped appropriately.			

```

topology disksim_iodriver driver0 [
  disksim_bus bus0 [
    disksim_ctrlr ctrlr0 [
      disksim_bus bus1 [
        disksim_disk disk0 []
      ] # end of bus1
      disksim_bus bus2 [
        disksim_disk disk1 []
        disksim_disk disk2 []
        disksim_disk disk3 []
      ] # end of bus2
    ] # end of ctrlr0
  ] # end of bus0
] # end of system topology

```

Figure 1: A Storage System Topology

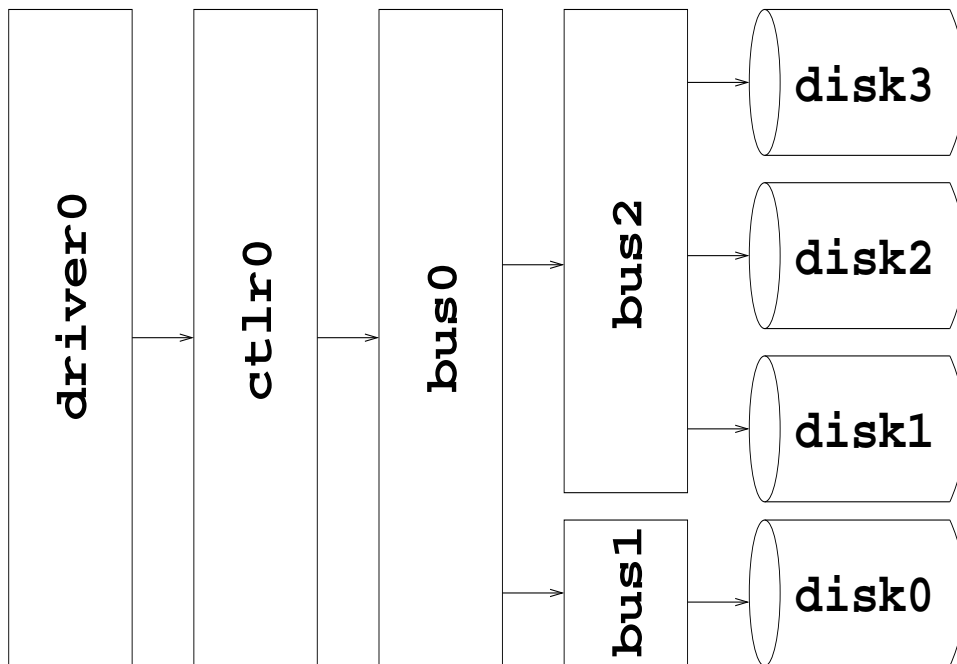


Figure 2: Diagram of the storage system corresponding to the above specification

<code>disksim_logorg</code>	Distribution scheme	string	required
This specifies the data distribution scheme (which is orthogonal to the redundancy scheme). <code>Asis</code> indicates that no re-mapping occurs. <code>Striped</code> indicates that data are striped over the organization members. <code>Random</code> indicates that a random disk is selected for each request. N.B.: This is only to be used with constant access-time disks for load-balancing experiments. <code>Ideal</code> indicates that an idealized data distribution (from a load balancing perspective) should be simulated by assigning requests to disks in a round-robin fashion. Note that the last two schemes do not model real data layouts. In particular, two requests to the same block will often be sent to different devices. However, these data distribution schemes are useful for investigating various load balancing techniques [3]. N.B.: This is only to be used with constant access-time disks for load-balancing experiments.			
<code>disksim_logorg</code>	Redundancy scheme	string	required
This specifies the redundancy scheme (which is orthogonal to the data distribution scheme). <code>Noredun</code> indicates that no redundancy is employed. <code>Shadowed</code> indicates that one or more replicas of each data disk are maintained. <code>Parity_disk</code> indicates that one parity disk is maintained to protect the data of the other organization members. <code>Parity_rotated</code> indicates that one disk's worth of data (spread out across all disks) are dedicated to holding parity information that protects the other N-1 disks' worth of data in an N-disk organization.			
<code>disksim_logorg</code>	Components	string	optional
This specifies whether the data organization's component members are entire disks ( <code>Whole</code> ) or partial disks ( <code>Partial</code> ). Only the former option is supported in the first released version of DiskSim.			
<code>disksim_logorg</code>	devices	list	required
List of device names to be included in this logical organization.			
<code>disksim_logorg</code>	Stripe unit	int	required
This specifies the stripe unit size. 0 indicates fine-grained striping (e.g., bit or byte striping), wherein all data disks in the logical organization contain an equal fraction of every addressable data unit.			
<code>disksim_logorg</code>	Synch writes for safety	int	required
This specifies whether or not an explicit effort should be made to do the N+1 writes of a parity-protected logical organization at "the same time" when handling a front-end write request with the read-modify-write (RMW) approach to parity computation. If true (1), then all reading of old values (for computing updated parity values) must be completed before the set of back-end writes is issued. If false (0), then each back-end write is issued immediately after the corresponding read completes (perhaps offering improved performance).			
<code>disksim_logorg</code>	Number of copies	int	required
This specifies the number of copies of each data disk if the logical organization employs <code>Shadowed</code> redundancy. Otherwise, this parameter is ignored.			
<code>disksim_logorg</code>	Copy choice on read	int	required
This specifies the policy used for selecting which disk from a set of <code>Shadowed</code> replicas should service a given read request since any of them can potentially do so. 1 indicates that all read requests are sent to a single primary replica. 2 indicates that one of the replicas should be randomly selected for each read request. 3 indicates that requests should be assigned to replicas in a round-robin fashion. 4 indicates that the replica that would incur the shortest seek distance should be selected and ties are broken by random selection. 5 indicates that the replica that has the shortest request queue should be selected and ties are broken by random selection. 6 indicates that the replica that has the shortest request queue should be selected and ties are broken by policy 4 (see above). This parameter is ignored if <code>Shadowed</code> replication is not chosen.			

disksim_logorg	RMW vs. reconstruct	float	required
This specifies the breakpoint in selecting Read-Modify-Write (RMW) parity updates (verses complete reconstruction) as the fraction of data disks that are updated. If the number of disks updated by the front-end write request is smaller than the breakpoint, then the RMW of the “old” data, “old” parity, and “new” data is used to compute the new parity. Otherwise, the unmodified data in the affected stripe are read from the corresponding data disks and combined with the new data to calculate the new parity. This parameter is ignored unless some form of parity-based replication is chosen.			

disksim_logorg	Parity stripe unit	int	required
This specifies the stripe unit size used for the <code>Parity_rotated</code> redundancy scheme. This parameter is ignored for other schemes. The parity stripe unit size does not have to be equal to the stripe unit size, but one must be a multiple of the other. Use of non-equal stripe unit sizes for data and parity has not been thoroughly tested in the current release of DiskSim.			

disksim_logorg	Parity rotation type	int	required
This specifies how parity is rotated among the disks of the logical organization. The four options, as described in [11], are 1 - left symmetric, 2 - left asymmetric, 3 - right asymmetric, 4 - right symmetric. This parameter is ignored unless <code>Parity_rotated</code> redundancy is chosen.			

disksim_logorg	Time stamp interval	float	required
This specifies the interval between “time stamps.” A value of 0.0 for this parameter disables the time stamp mechanism.			

disksim_logorg	Time stamp start time	float	required
This specifies the simulated time (relative to the beginning of the simulation) of the first time stamp.			

disksim_logorg	Time stamp stop time	float	required
This specifies the simulated time (relative to the beginning of the simulation) of the last time stamp.			

disksim_logorg	Time stamp file name	string	required
This specifies the name of the output file to contain a log of the instantaneous queue lengths of each of the organization’s back-end devices at each time stamp. Each line of the output file corresponds to a single time stamp and contains the queue lengths of each device separated by white space. A value of “0” or of “null” disables this feature (as does disabling the time stamp mechanism).			

The “time-stamp” parameters configure DiskSim’s per-logorg mechanism for collecting information about instantaneous per-device queue lengths at regular intervals.

### 3.9 Process-Flow Parameters

The various parameters involved with configuring the synthetic workload generation module are described in Section 4.2.

disksim_pf	Number of processors	int	required
This specifies the number of processors used by the simple system-level model. These processors (and, more generally, DiskSim’s system-level model) are only used for the synthetic generation module.			

disksim_pf	Process-Flow Time Scale	float	required
This specifies a multiplicative scaling factor for computation times “executed” by a simulated processor. For example, 2.0 doubles each computation time, and 0.5 halves each computation time.			



## 4 Input Workloads: Traces and Synthetic Workloads

DiskSim can be exercised with I/O requests in several ways, including external traces, internally-generated synthetic workloads, and interactions within a larger simulation environment (e.g., a full system simulator). This section describes each of these options.

### 4.1 Traces

DiskSim can accept traces in several formats, and new formats can be added with little difficulty. This subsection describes the default input format and briefly describes how to add support for new trace formats. The DiskSim 1.0 distribution supports the default format (“ascii”), a validation trace format (“validate”), the raw format (“raw”) of the disk request traces described in [5, 3], and the raw format (“hpl”, or “hpl2” if the trace file header has been stripped) of the disk request traces described in [15].

#### 4.1.1 Default Format

The default input format is a simple ASCII stream (or file), where each line contains values for five parameters (separated by white space) describing a single disk request. The five parameters are:

1. **Request arrival time:** Float [nonnegative milliseconds] specifying the time the request “arrives” relative to the start of the simulation (at time 0.0). Requests must appear in the input stream in ascending time order.
2. **Device number:** Integer specifying the device number (i.e., the storage component that the request accesses). The device mappings (see Section 3), if any, are applied to this value.
3. **Block number:** Integer [nonnegative] specifying the first device address of the request. The value is specified in the appropriate access unit of the logical device in question, which may be modified by the device mappings (see Section 3).
4. **Request size:** Integer [positive] specifying the size of the request in device blocks (i.e., the access unit of the logical device in question).
5. **Request flags:** Hex Integer comprising a Boolean Bitfield specifying additional information about the request. For example, bit 0 indicates whether the request is a read (1) or a write (0). Other bits specify information that is most appropriate to a full system simulation environment (e.g., request priority). Valid bitfield values are listed in “disksim\_global.h”.

An example trace in this format is included with the distribution.

#### 4.1.2 Adding Support For New Trace Formats

Adding support for a new trace format requires only a few steps:

1. Add a new trace format constant to “disksim\_global.h”.
2. Select a character string to represent the trace format on the command line. Add a format name comparison to “iotrace\_set\_format” in “disksim\_iotrace.c”.
3. Create a procedure (“iotrace\_XXXX\_get\_ioreq\_event”) in “disksim\_iotrace.c” to read a single disk request description from an input trace of the new format (“XXXX”) and construct a disk request event in the internal format (described briefly below). The functions “iotrace\_read\_[char,short,int32]” can simplify this process. Incorporate the new function into main switch statement in “iotrace\_get\_ioreq\_event”. The internal DiskSim request structure (at request arrival time) is not much more complicated than the default (ascii) trace format. It contains “time,” “devno,” “blkno,” “bcount,” and “flags” fields that correspond to the five non-auxiliary fields described above. The other fields do not have to be initialized, except that “opid” should be zeroed. See “iotrace\_ascii\_get\_ioreq\_event” for an example.
4. If the trace file has an informational header (useful or otherwise), then create a procedure (“iotrace\_XXXX\_initialize\_file”) in “disksim\_iotrace.c” and add it into the if/else statement in “iotrace\_initialize\_file”.

## 4.2 Synthetic Workloads

DiskSim includes a simple synthetic workload generating module that can be used to drive storage subsystem simulations. The parameter file specifies the number of generators (similar to processes) and the characteristics of the workloads generated by each. Each synthetic generator “executes” as a process in a very simple system-level model, issuing I/O requests after periods of “think time” and, when appropriate, waiting for them to complete. This module can be configured to generate a wide range of synthetic workloads, both with respect to the disk locations accessed and the request arrival times.

### 4.2.1 Configuration

The synthetic generation module is configured with the parameters (described below) specified in the last section of the parameter file. The other “Process-Flow Input Parameters” are also relevant. In particular, the number and time scaling of the simulated processors is important, since these processors “execute” the inter-request think times as computation times (one process’s computation per processor at a time).

<code>disksim_synthio</code>	Number of I/O requests to generate	int	required
This specifies the number of independent, concurrent, request-generating processes.			

<code>disksim_synthio</code>	Maximum time of trace generated	float	required
This specifies the maximum number of I/O requests to generate before ending the simulation run. A simulation run continues until either the specified number of requests is generated or the maximum simulation time (see below) is reached.			

<code>disksim_synthio</code>	System call/return with each request	int	required
This specifies whether or not each request occurs within the context of a system call (which may affect the behavior of the associated process in the system-level model). If true (1), each request will be preceded by a system call event and followed by a system call return event.			

<code>disksim_synthio</code>	Think time from call to request	float	required
This specifies the think time (i.e., computation time) between the system call event and the disk request event. This parameter is only relevant if the above Boolean parameter is set to true (1).			

<code>disksim_synthio</code>	Think time from request to return	float	required
This specifies the think time (i.e., computation time) between the disk request event and the system call return event. This parameter is only relevant if the above Boolean parameter is set to true (1).			

<code>disksim_synthio</code>	Generators	list	required
A list of synthgen block values describing the generators.			

The generators (described below) contain several probability-distribution parameters which are defined as follows: A distribution is defined by a list whose first element is a string naming the type of distribution; the subsequent list elements are specific to the type of distribution.

The distribution name may be one of:

- `uniform` requiring two additional floats specifying minimum and maximum values
- `normal` requiring two additional Floats specifying mean and variance values. As the second value is a variance, it must be nonnegative.
- `exponential` requiring two additional Floats specifying base and mean values.
- `poisson` requiring two additional Floats specifying base and mean values.

- `twovalue` requiring three additional Floats specifying a default value, a secondary value, and a probability indicating how often the secondary value should be returned. As the last Float is a measure of probability, it must be between 0.0 and 1.0.

All of the distributions are computed as specified, but value generation is repeated whenever an illegal value results (e.g., a negative inter-arrival time).

<code>disksim.synthgen</code>	Storage capacity per device	int	required
This specifies the number of unique storage addresses per storage device (in the corresponding device's unit of access) accessible to generators of this type.			

<code>disksim.synthgen</code>	<code>devices</code>	list	required
This specifies the set of storage devices accessible to generators of this type. The devices may be either the names of individual devices in a "parts" logorg or the name of an "array" logorg.			

<code>disksim.synthgen</code>	Blocking factor	int	required
This specifies a unit of access for generated requests that is a multiple of the storage devices' unit of access. All generated request starting addresses and sizes will be a multiple of this value.			

<code>disksim.synthgen</code>	Probability of sequential access	float	required
This specifies the probability that a generated request is sequential to the immediately previous request. A sequential request starts at the address immediately following the last address accessed by the previously generated request.			

<code>disksim.synthgen</code>	Probability of local access	float	required
This specifies the probability that a generated request is "local" to the immediately previous request. A local request begins some short distance away from the previous request's starting address, where the distance is computed via a random variable definition described below.			

<code>disksim.synthgen</code>	Probability of read access	float	required
This specifies the probability that a generated request is a read.			

<code>disksim.synthgen</code>	Probability of time-critical request	float	required
This specifies the probability that a generated request is time-critical. That is, the corresponding generator process "blocks" and waits for the request to complete before continuing with its sequence of work (i.e., its next think time) [5, 2].			

<code>disksim.synthgen</code>	Probability of time-limited request	float	required
This specifies the probability that a generated request is time-limited. That is, the corresponding generator process "blocks" and waits for the request to complete (if it is not already complete) after a given amount of think time (specified by the below "time limit" parameters) [5, 2].			

<code>disksim.synthgen</code>	Time-limited think times	list	required
This is a random variable distribution specifying the time limit for a time-limited request. Note that the generated time limit (i.e., the computation time occurring before the generator process "blocks" and waits for the request to complete) may differ from the actual time limit (due to CPU contention).			

<code>disksim.synthgen</code>	General inter-arrival times	list	required
This is a random variable distribution specifying the inter-request think time preceding the generated request if the generated request's starting address is unrelated to the previous request's starting address (i.e., if the generated request's address is "random" rather than "sequential" or "local").			

<code>disksim_synthgen</code>	Sequential inter-arrival times	list	required
This is a random variable distribution specifying the inter-request think time preceding the generated request if the generated request's starting address is "sequential" to the previous request's starting address.			
<code>disksim_synthgen</code>	Local inter-arrival times	list	required
This is a random variable distribution specifying the inter-request think time preceding the generated request if the generated request's starting address is "local" to the previous request's starting address.			
<code>disksim_synthgen</code>	Local distances	list	required
This is a random variable distribution specifying the distance from the previous request's starting address when generating a "local" request's starting address.			
<code>disksim_synthgen</code>	Sizes	list	required
This is a random variable distribution specifying the request size.			

Probability of sequential access along with Probability of local access determine how a generated request's starting address is assigned. Their sum must be less than or equal to 1.0. Each request's starting address is sequential, local or random. A random request is assigned a device and starting address from a uniform distribution spanning the entire available storage space (as specified by the above parameters).

Probability of time-critical request along with Probability of time-limited request determine how a generated request's criticality is assigned. Their sum must be less than or equal to 1.0. Each request is time-critical, time-limited or time-noncritical. A generator process never "blocks" and waits for time-noncritical requests; it simply generates them and continues with its work.

With these probabilities, the synthetic generators can be configured to emulate open and closed subsystem models, as well as a range of intermediate options. An open subsystem model (like most trace-driven simulations) generates new requests independently of the completion times of previous requests. It can be created by setting the probabilities of time-critical and time-limited requests both to zero (0.0) and configuring the system-level model to have the same number of processors as there are generators. The inter-request computation times will therefore be the inter-arrival times (per generator). A closed subsystem model (like many queuing models) generates a new request only after the previous request completes, keeping a constant number of requests in the system (either hidden in think times or being serviced). It can be created by setting the probability that a request is time-critical to one (1.0). Setting the inter-request computation times (below) to zero (0.0) eliminates think times and results in a constant number of requests "in service" at all times.

### 4.3 Incorporating DiskSim Into System-Level Simulators

With a modicum of effort, DiskSim can be incorporated into a full system-level simulator in order to provide accurate timings for the handling of storage I/O requests. This section briefly describes one method for doing so, and "disksim\_interface.c" stubs out a possible implementation. Using this approach (which assumes only a few characteristics of the system simulator), DiskSim will act as a slave of the system simulator, providing disk request completion indications in time for an interrupt to be generated in the system simulation. Specifically, DiskSim code will only be executed when invoked by one of the following procedures, called as appropriate by the containing simulator:

- **disksim\_initialize**: for initializing the DiskSim state.
- **disksim\_shutdown**: for printing statistics at the end of a simulation run.
- **disksim\_dump\_stats**: for printing the current running statistics during the simulation (e.g., at a statistics check-point).
- **disksim\_internal\_event**: for "calling back" into DiskSim so that it can update its internal simulation "time" to match the system-level simulator's global "time" and handle any internal DiskSim events that occur in the intervening time. Note that this function is called by the system-level simulator *on behalf* of DiskSim, since DiskSim no longer has control over the global simulation "time." Additional details are given below.

- **disksim\_request\_arrive**: for issuing an I/O request into DiskSim.

Using this interface requires only two significant functionalities of the system simulation environment:

1. The ability to function correctly without knowing when a disk request will complete at the time that it is initiated. The system simulation will be informed at some later point (in its view of time) that the request is completed. At this time, the appropriate “disk request completion” interrupt could be inserted into the system simulation.
2. The ability for DiskSim to register callbacks with the system simulation environment. That is, this interface code must be able to request (of the system-level simulator) an invocation of a callback function (such as `disksim_internal_event`, described above) when the simulated time reaches a DiskSim-specified value. It is also helpful (but not absolutely necessary) to be able to “de-schedule” a callback at some point after it has been requested. For example, a callback requested to indicate the end of some disk prefetching activity may be superceded by a new request arriving at the disk (and interrupting the ongoing prefetch).

If the actual content on the disk media (i.e., the “data”) must be maintained during the course of a system-level simulation, this functionality can easily be provided by code outside of DiskSim, which does not itself provide such functionality.

## 5 The Output File

At the beginning of the simulation, the values for the numerous configuration parameters are copied into the beginning of the output file. The remainder of the output file contains the aggregate statistics of the simulation run, including both the characteristics of the simulated workload (if enabled) and the performance indicators of the various storage components. Each line of the latter portion of the output file (excluding the statistic distributions) is unique, simplifying the process of searching through the file for a particular result.

DiskSim collects a large number of statistics about the simulated storage components. As discussed in section 3, the size of the output file can be reduced by configuring DiskSim not to report undesired sets of statistics. As each statistic is described below, it will also be noted whether or not the statistic can be pruned from the output file via corresponding Boolean (enable/disable) input parameters.

The following subsection describes the configuration of statistics-collection parameters. The remainder of the section describes the sets of statistics which DiskSim collects.

### 5.1 The `statdefs` File

Although some of the results collected by DiskSim are simple counts or sums, aggregate statistics (average, standard deviation, distribution) are collected for many values. In particular, statistics reported as distributions make use of the `statdefs` file for configuring the “bins” which count observed values falling within specific ranges. Each statistic description consists of four lines, as described below. Comments or other extraneous information may be placed within the file as long as each 4-line description is contiguous. Most of the statistics in the sample `statdefs` file provided in the DiskSim distribution are mandatory; DiskSim will emit an error message and exit if any such statistic is not defined. Each statistic definition consists of:

```
<statname>
Distribution size: <distsize>
Scale/Equals: <scale>/<equals>
<bins>
```

`statname` is the name of the statistic being described. DiskSim searches for this value on a line by itself to identify the beginning of the corresponding statistic description. So, the order of the statistics in the `statdefs` file is unimportant.

`distsize` is a positive integer specifying the number of bins into which to partition the observed values. When this value is less than or equal to the internal DiskSim constant `DISTSIZE`, the entire distribution is reported on two lines (for output file compactness and readability). The output for larger distributions consists of one bin per line.

`scale` is a nonzero multiplication factor to be applied to observed values before selecting an appropriate bin. It is useful mainly because the bin boundaries are specified as integers (see below). For example, if a specific “response time” statistic’s distribution is specified using microsecond-based boundaries, the “Scale” value should be set to 1000 (since internal times in DiskSim are millisecond-based).

`equals` is a nonnegative integer that specifies how many of the bins (starting from the first one) should collect only observed values that are exactly equal to the specified boundary, rather than values that are less than the boundary (see below).

`bins` describes the bin boundaries in one of two formats. If the number of bins is less than or equal to the internal DiskSim constant `DISTSIZE`, then this line contains nine integers. If the `equals` value (see above) is set to some value  $n$ , the first  $n$  integers specify bins that will hold exact observed values rather than ranges of observed values. The remaining integers specify bins holding observed values below the specified integer boundary. When DiskSim categorizes an observed value, the bins are checked in the sequence in which they are found on the line. So, bin boundaries beyond the “Equals” values should be in ascending order to avoid unclear results.

If the number of bins is greater than `DISTSIZE`, the format of this line is a string

```
Start <x> step <y> grow <z>
```

where  $x$  specifies the first bin boundary,  $y$  specifies a constant value to be added to the current boundary value when computing the next boundary value, and  $z$  specifies a percent of the previous bin value to be added to the current boundary value when computing the next boundary value. The combination of a multiplicative scaling factor (“Scale”),

an additive step function (“step”), and a multiplicative step function (“grow”) provides the flexibility to concisely specify a wide range of typical statistic distributions.

## 5.2 Simulation Results

Each statistic in the output file is identified by a unique string. A typical single-value statistic is reported as “<name> : <value>.” In some cases a single value is inadequate to describe the statistic, so a set of four aggregate statistics are reported instead: the average value, the standard deviation, the maximum value, and the distribution as specified by the corresponding entry in the statdef file.

The format for reporting the first three aggregates mentioned above is

```
<context> <statname> <type>: <value>
```

where the `context` is the context of the statistic (e.g., “Disk #2”), `statname` is the name of the statistic used to index into the statdefs file, and `type` is the aggregate type [`average`, `std.dev.`, or `maximum`]. The corresponding result value is specified on the same line. The format for reporting distributions is similar to the other three aggregates, but the value is reported on lines following the indentifying string. If the number of bins specified in the corresponding entry in the statdefs file is greater than `DISTSIZE`, then the distribution is reported on two lines: one for the bin boundaries and one for the observed bin counts. Otherwise, the distribution is reported on the next  $n$  lines, where  $n$  is the number of bins. Each line contains four values: (1) the bin boundary, (2) the count of observed values falling into the bin, (3) the measured probability that an observed value falls into the bin (i.e., the count divided by the number of observations), and (4) the measured probability that an observed value falls in the bin or any previous bin.

Each individual statistic found in the output file is described below. *Unless otherwise specified, all statistics measuring simulation “time” are reported in milliseconds.*

**Total time of run:** The simulated time at the end of the simulation run.

**Warm-up time:** The simulation warm-up time not covered by the reported statistics. See Section 3.1.

The next set of statistics falls into one of several categories, depending on the style of input. If synthetic generation is used, then statistics for the simple system-level model are reported (Section 5.2.1). If a validation trace is used, then statistics about the behavior measured for the real disk are reported (Section 5.2.2). If the traces described in [15] (referred to in DiskSim as “HPL” traces) are used, then statistics about performance observed for the traced system are reported (Section 5.2.3). If any other external trace is used, then no statistics are reported in this section of the output file, and the I/O driver statistics are next in the output file (Section 5.2.5), followed by the disk drive statistics (Section 5.2.6), controller statistics (Section 5.2.7), and bus statistics (Section 5.2.8).

### 5.2.1 Process-flow Statistics

The following statistics are reported only when the internal synthetic workload generator is enabled.

**CPU Total idle milliseconds:** the sum of the idle times for all simulated CPUs in the process-flow simulation.

**CPU Idle time per processor:** the average per-CPU idle time.

**CPU Percentage idle cycles:** the average percentage of time that each CPU spent idle.

**CPU Total false idle ms:** the sum of the false idle times for all CPUs. “False idle time” is that time that a processor spends idle because processes are blocked waiting for I/O (e.g., disk requests), as opposed to real idle time during which there is no work for the CPU to do.

**CPU Percentage false idle cycles:** the average percentage of each CPUs time that was consumed by false idle time.

**CPU Total idle work ms:** the sum of the idle work times for all CPUs. “Idle work time” is useful computation (e.g., interrupt handlers and background tasks) completed while in the idle state (because no user processes are runnable).

**CPU Context Switches:** the number of context switches.

**CPU Time spent context switching:** the aggregate CPU time consumed by context switch overheads on all CPUs.

**CPU Percentage switching cycles:** the average percentage of each CPUs time that was consumed by context switching overheads.

**CPU Number of interrupts:** the total number of interrupts received by all CPUs.

**CPU Total time in interrupts:** the total computation time consumed by interrupt handlers.

**CPU Percentage interrupt cycles:** the average percentage of each CPU's time that was consumed by interrupt handling.

**CPU Time-Critical request count:** the total number of time-critical requests generated.

**CPU Time-Critical Response time stats:** aggregate statistics for the response times observed for all time-critical requests.

**CPU Time-Limited request count:** the total number of time-limited requests generated.

**CPU Time-Limited Response time stats:** aggregate statistics for the response times observed for all time-limited requests.

**CPU Time-Noncritical request count:** the total number of time-noncritical requests generated.

**CPU Time-Noncritical Response time stats:** aggregate statistics for the response times observed for all time-noncritical requests.

The next four statistics are not reported if "Print all interrupt stats?" is set to false (0).

**CPU Number of IO interrupts:** the total number of I/O interrupts received by all CPUs.

**CPU Time spent in I/O interrupts:** the total computation time spent on I/O interrupt handlers.

**CPU Number of clock interrupts:** the total number of I/O interrupts received by all CPUs.

**CPU Time spent in clock interrupts:** the total computation time spent on clock interrupt handlers.

The next four statistics are not reported if "Print sleep stats?" is set to false (0).

**Number of sleep events:** the total number of sleep events "executed" by all processes.

**Number of I/O sleep events:** the total number of sleep events "executed" in order to wait for I/O requests.

**Average sleep time:** the average length of time between a sleep event and the corresponding wake-up event.

**Average I/O sleep time:** the average length of time between a sleep event that waits for an I/O request and the corresponding wake-up event (i.e., the I/O request's completion).

If there is more than one CPU, then per-CPU statistics are reported. The per-CPU statistics are the same as the aggregate CPU statistics described above. The per-CPU statistics are not reported if "Print per-CPU stats?" is set to false (0).

**Process Total computation time:** the total computation time of all processes (other than the idle processes) in the system.

**Process Last event time:** the simulated time of the last process event "executed."

**Process Number of I/O requests:** the total number of I/O requests generated.

**Process Number of read requests:** the total number of read I/O requests generated.

**Process Number of C-switches:** the total number of context switches to or from non-idle processes.

**Process Number of sleeps:** the total number of sleep events "executed."

**Process Average sleep time:** the average time between a process's sleep event and the corresponding wake-up event.

**Process Number of I/O sleeps:** the total number of sleep events "executed" in order to wait for I/O requests.

**Process Average I/O sleep time:** the average time between a process's sleep event that waits for an I/O request and the corresponding wake-up event (i.e., the I/O request's completion).

**Process False idle time:** the total amount of false idle time. This value can be greater than the total measured false idle time because more than one process can contribute to any given period of false idle time.

**Process Read Time limits measured:** the number of time limits observed for read I/O requests. (This value differs from the number of time-limited read requests if the simulation ends before one or more read I/O request time limits expire.)

**Process Read Time limit duration stats:** aggregate statistics for read I/O request time limits.

**Process Write Time limits measured:** the number of time limits measured for write I/O requests. (This value differs from the number of time-limited write requests if the simulation ends before one or more write I/O request time limits expire.)

**Process Write Time limit duration stats:** aggregate statistics for write I/O request time limits.

**Process Read Time limits missed:** the number of time limits missed by read I/O requests.

**Process Missed Read Time limit duration stats:** aggregate statistics for missed read I/O request time limits.

**Process Write Time limits missed:** the number of time limits missed by write I/O requests.

**Process Missed Write Time limit duration stats:** aggregate statistics for the missed write I/O request time limits.



If there is more than one simulated process, then per-process statistics are reported. The per-process statistics are the same as the aggregate process statistics described above. The per-process statistics are not reported if “Print per-process stats?” is set to false (0).

### 5.2.2 Validation Trace Statistics

The following statistics are reported only when an external validation trace (with a format of “validate”) is used as the input workload.

**VALIDATE Trace access time stats:** aggregate statistics for the access times measured for the corresponding real storage subsystem. (The access time measured for each request is part of the input trace format.)

**VALIDATE Trace access diff time stats:** aggregate statistics for the per-request differences between the simulated and measured access times.

**VALIDATE Trace write access diff time stats:** aggregate statistics for the measured access times for write requests.

**VALIDATE Trace write access diff time stats:** aggregate statistics for the per-request differences between the simulated and measured access times for write requests.

The remaining statistics for validate workloads are historical in nature and are primarily useful for debugging DiskSim’s behavior. The information needed to trigger them is not included in most of the validation traces.

**VALIDATE double disconnects:** the number of requests incurring two bus disconnects during the request’s lifetime.

**VALIDATE triple disconnects:** the number of requests incurring three bus disconnects during the request’s lifetime.

**VALIDATE read buffer hits:** the number of read requests that were serviced directly from the disk’s on-board cache.

**VALIDATE buffer misses:** the number of requests that required actual magnetic media access.

### 5.2.3 HPL Trace Statistics

The following statistics are reported only when an external HPL trace (i.e., a trace in the HPLabs SRT format) is used as the input workload.

**Total reads:** the number of read requests in the trace, followed by the fraction of all requests that were reads.

**Total writes:** the number of write requests in the trace, followed by the fraction of all requests that were writes.

**Sync Reads:** the number of read requests marked (by a flag value) as synchronous, meaning that an application process will be blocked until the request completes. This value is followed by the fraction of all requests that were synchronous reads and the fraction of all read requests that were synchronous.

**Sync Writes:** the number of write requests marked (by a flag value) as synchronous, meaning that an application process will be blocked until the request completes. This value is followed by the fraction of all requests that were synchronous writes and the fraction of write requests that were synchronous.

**Async Reads:** the number of read requests not marked as synchronous, followed by the fraction of requests that were asynchronous reads and the fraction of all read requests that were asynchronous.

**Async Writes:** the number of write requests not marked as synchronous, followed by the fraction of all requests that were asynchronous writes and the fraction of all write requests that were asynchronous.

**Mapped disk #X Trace queue time stats:** aggregate statistics for the per-request queue times measured for disk X in the traced system.

**Mapped disk #X Trace response time stats:** aggregate statistics for the per-request response times (i.e., request arrival to request complete, including queue delays and service time) measured for disk X in the traced system.

**Mapped disk #X Trace access time stats:** aggregate statistics for the per-request access times (i.e., service times) measured for disk X in the traced system.

**Mapped disk #X Trace queue length stats:** aggregate statistics for the instantaneous queue lengths observed by each request for disk X in the traced system.

**Mapped disk #X Trace non-queue time stats:** aggregate statistics for the measured per-request times spent in the device driver for disk X in the traced system between request arrival and delivery to the storage controller when no other requests are pending. This provides insight into the device driver overheads for the traced systems.

### 5.2.4 System-level Logical Organization Statistics

The following statistics are reported for each system-level logical organization. (Note: Every DiskSim simulation involves at least one system-level logical organization, even if it simply maps each logical device onto an equivalent physical device.)

**System logorg #X Number of requests:** the number of requests submitted to the front-end of logical organization x.

**System logorg #X Number of read requests:** the number of read requests submitted to the front-end of logical organization X, followed by the fraction of all front-end requests that were reads.

**System logorg #X Number of accesses:** the number of accesses passed to the back-end of logical organization X. Striping, data redundancy, and other logical aspects can cause this value to differ from the number of front-end requests.

**System logorg #X Number of read accesses:** the number of accesses passed to the back-end of logical organization X in response to front-end read requests, followed by the fraction of all back-end requests that were reads.

**System logorg #X Average outstanding:** the average number of front-end requests in progress at any point in time to logical organization X.

**System logorg #X Maximum outstanding:** the maximum number of front-end requests in progress at any point in time to logical organization X.

**System logorg #X Avg nonzero outstanding:** the average number of front-end requests in progress at times when there was at least one outstanding request to logical organization X.

**System logorg #X Completely idle time:** the amount of time during which no front-end requests are outstanding to logical organization X.

**System logorg #X Response time stats:** aggregate statistics for the response times observed for all front-end requests to logical organization X.

**System logorg #X Time-critical reads:** the number of front-end read requests to logical organization X marked (by a flag field) as time-critical.

**System logorg #X Time-critical write:** the number of front-end write requests to logical organization X marked (by a flag field) as time-critical.

The next ten statistics are not reported if “Print driver locality stats?” is set to false (0).

**System logorg #X Inter-request distance stats:** aggregate statistics for the distances between the starting addresses of subsequent accesses to the same device in logical organization X.

**System logorg #X Sequential reads:** the number of back-end read accesses whose starting addresses were sequential to the immediately previous access to the same device in logical organization X, followed by the fraction of back-end accesses that were sequential reads and the fraction of back-end reads that were sequential.

**System logorg #X Sequential writes:** the number of back-end write accesses whose starting addresses were sequential to the immediately previous access to the same device in logical organization X, followed by the fraction of back-end accesses that were sequential writes and the fraction of back-end writes that were sequential.

**System logorg #X Interleaved reads:** the number of back-end read accesses whose starting addresses were almost sequential to (i.e., less than 16 sectors beyond the end of) the immediately previous access to logical organization X, followed by the fraction of back-end accesses that were “interleaved” reads and the fraction of back-end reads that were “interleaved.”

**System logorg #X Interleaved writes:** the number of back-end write accesses whose starting addresses were almost sequential to (i.e., less than 16 sectors beyond the end of) the immediately previous access to logical organization X, followed by the fraction of back-end accesses that were “interleaved” writes and the fraction of back-end writes that were “interleaved.”

**System logorg #X Logical sequential reads:** the number of front-end read requests whose starting addresses were logically sequential to the immediately previous request to logical organization X.

**System logorg #X Logical sequential writes:** the number of front-end write requests whose starting addresses were logically sequential to the immediately previous request to logical organization X.

**System logorg #X Sequential disk switches:** the number of back-end accesses generated for logically sequential front-end requests for logical organization X that (because of striping or some such) accessed a different device than the immediately previous request.

**System logorg #X Logical local accesses:** the number of front-end requests marked (by a flag) as logically “local” to the immediately previous request to logical organization X.

**System logorg #X Local disk switches:** the number of back-end accesses generated for front-end requests marked (by a flag) as logically “local” that (because of striping or some such) accessed a different device than the immediately previous request to logical organization X.

The next two statistics are not reported if “Print driver interfere stats?” is set to false (0).

**System logorg #X Sequential step S:** the number of back-end accesses to logical organization X that were sequential to the back-end access S+1 accesses prior, followed by the fraction of back-end accesses that fall into this category. These statistics are only reported if the fraction is greater than 0.002.

**System logorg #X Local (D) step S:** the number of back-end accesses to logical organization X whose starting addresses were D device sectors from the back-end access S+1 accesses prior, followed by the fraction of back-end accesses that fall into this category. These statistics are only reported if the fraction is greater than 0.002.

The next two statistics are not reported if “Print driver blocking stats?” is set to false (0).

**System logorg #X Blocking factor: B:** the number of back-end accesses to logical organization X whose size is an integer multiple of B sectors, followed by the fraction of back-end accesses that fall into this category. These statistics are only reported if the fraction is greater than 0.002.

**System logorg #X Alignment factor: A:** the number of back-end accesses to logical organization X whose starting address is an integer multiple of A sectors, followed by the fraction of back-end accesses that fall into this category. These statistics are only reported if the fraction is greater than 0.002.

The next three statistics are not reported if “Print driver intarr stats?” is set to false (0).

**System logorg #X Inter-arrival time stats:** aggregate statistics for the inter-arrival times of front-end requests to logical organization X.

**System logorg #X Read inter-arrival time stats:** aggregate statistics for the inter-arrival times of front-end read requests to logical organization X.

**System logorg #X Write inter-arrival time stats:** aggregate statistics for the inter-arrival times of front-end write requests to logical organization X.

The next two statistics are not reported if “Print driver streak stats?” is set to false (0).

**System logorg #X Number of streaks:** the number of sequences of back-end accesses to logical organization X addressed to the same device with no interleaved accesses to other devices.

**System logorg #X Streak length stats:** aggregate statistics for the lengths of sequences of back-end accesses to logical organization X addressed to the same device with no interleaved accesses to other devices.

The next three statistics are not reported if “Print driver stamp stats?” is set to false (0).

**System logorg #X Timestamped # outstanding distribution:** a distribution of the number of requests outstanding to logical organization X at regular simulated time intervals (specified by the “Time stamp interval” parameter).

**System logorg #X Timestamped avg # outstanding difference distribution:** a distribution of the average difference between the number of requests outstanding to each back-end device of logical organization X and the average number of requests outstanding per back-end device of logical organization X (measured at regular simulated time intervals).

**System logorg #X Timestamped max # outstanding difference distribution:** a distribution of the maximum difference between the number of requests outstanding to a particular back-end device of logical organization X and the average number of requests outstanding per back-end device of logical organization X (measured at regular simulated time intervals).

The next three statistics are not reported if “Print driver size stats?” is set to false (0).

**System logorg #X Request size stats:** aggregate statistics for the sizes of front-end requests to logical organization X.

**System logorg #X Read request size stats:** aggregate statistics for the sizes of front-end read requests to logical organization X.

**System logorg #X Write request size stats:** aggregate statistics for the sizes of front-end write requests to logical organization X.

The next two statistics are not reported if “Print driver idle stats?” is set to false (0).

**System logorg #X Number of idle periods:** the number of time periods during which no requests were outstanding to logical organization X.

**System logorg #X Idle period length stats:** aggregate statistics for the durations of time periods during which no requests were outstanding to logical organization X.

The remaining system-level logorg statistics are aggregate statistics over the set of disks in the logical organization. The statistics reported are the same as those described in Section 5.2.6 under “Disk statistics.”

### 5.2.5 I/O Driver Statistics

All of the I/O driver statistics are generated by the request queue module (which is also used by the disk and controller modules). None of them are reported if “Print driver queue stats?” is set to false (0).

**IOdriver Total Requests handled:** the number of requests completed from the driver’s point of view.

**IOdriver Requests per second:** the number of requests completed per second of simulated time.

**IOdriver Completely idle time:** the total amount of time that no requests were outstanding.

**IOdriver Response time stats:** aggregate statistics for request response times.

**IOdriver Overlaps combined:** the number of requests made unnecessary because they completely overlap with another outstanding request, followed by the fraction of requests that fall into this category. (Note that this is an extremely unusual event in real systems, but the situation may arise frequently in trace-driven simulation [2].)

**IOdriver Read overlaps combined:** the number of read requests made unnecessary because they completely overlap with another outstanding request, followed by the fraction of requests that fall into this category.

The next eight statistics are not reported if “Print driver crit stats?” is set to false (0).

**IOdriver Critical Reads:** the number of read requests marked (by a flag) as time-critical, followed by the fraction of requests that are time-critical reads.

**IOdriver Critical Read Response time stats:** aggregate statistics for the response times of read requests marked time-critical.

**IOdriver Non-Critical Reads:** the number of read requests not marked (by a flag) as time-critical, followed by the fraction of requests that are reads not marked time-critical.

**IOdriver Non-Critical Read Response time stats:** aggregate statistics for the response times of read requests not marked time-critical.

**IOdriver Critical Writes:** the number of write requests marked (by a flag) as time-critical, followed by the fraction of requests that are time-critical writes.

**IOdriver Critical Write Response time stats:** aggregate statistics for the response times of write requests marked time-critical.

**IOdriver Non-Critical Writes:** the number of write requests not marked (by a flag) as time-critical, followed by the fraction of requests that are writes not marked time-critical.

**IOdriver Non-Critical Write Response time stats:** aggregate statistics for the response times of write requests not marked time-critical.

**IOdriver Number of reads:** the number of read requests, followed by the fraction of requests that are reads.

**IOdriver Number of writes:** the number of write requests, followed by the fraction of requests that are writes.

**IOdriver Sequential reads:** the number of read requests whose starting addresses are sequential to the immediately previous request to the same device, followed by the fraction of requests that are sequential reads.

**IOdriver Sequential writes:** the number of write requests whose starting addresses are sequential to the immediately previous request to the same device, followed by the fraction of requests that are sequential writes.

The next twelve statistics are not reported if “Print driver queue stats?” is set to false (0).

**IOdriver Average # requests:** the average number of requests outstanding (in queues or in service).

**IOdriver Maximum # requests:** the maximum number of requests outstanding.

**IOdriver end # requests:** the number of requests outstanding when the simulation ended.

**IOdriver Average queue length:** the average length of the request queue.

**IOdriver Maximum queue length:** the maximum length of the request queue.

**IOdriver End queued requests:** the length of the request queue when the simulation ended.

**IOdriver Queue time stats:** aggregate statistics for the queue times incurred by requests.

**IOdriver Avg # read requests:** the average number of read requests outstanding.

**IOdriver Max # read requests:** the maximum number of read requests outstanding.

**IOdriver Avg # write requests:** the average number of write requests outstanding.

**IOdriver Max # write requests:** the maximum number of write requests outstanding.

**IOdriver Physical access time stats:** aggregate statistics for the request access times (i.e., excluding any queuing times).

The next three statistics are not reported if “Print driver intarr stats?” is set to false (0).

**IOdriver Inter-arrival time stats:** aggregate statistics for request inter-arrival times.

**IOdriver Read inter-arrival time stats:** aggregate statistics for read request inter-arrival times.

**IOdriver Write inter-arrival time stats:** aggregate statistics for write request inter-arrival times.

The next two statistics are not reported if “Print driver idle stats?” is set to false (0).

**IOdriver Number of idle periods:** the number of time periods during which no requests were outstanding.

**IOdriver Idle period length stats** aggregate statistics for the durations of time periods during which no requests were outstanding.

The next three statistics are not reported if “Print driver size stats?” is set to false (0).

**IOdriver Request size stats:** aggregate statistics for the sizes of requests.

**IOdriver Read request size stats:** aggregate statistics for the sizes of read requests.

**IOdriver Write request size stats:** aggregate statistics for the sizes of write requests.

**IOdriver Instantaneous queue length stats:** aggregate statistics for the queue lengths observed at the points in time when each new request arrived.

**IOdriver Sub-optimal mapping penalty stats:** aggregate statistics for the seek distance penalties incurred due to the use of inaccurate mapping information in translating request starting locations to cylinder/track/sector locations (by a scheduler).

Some of the disk request scheduling algorithms supported by DiskSim employ more than one sub-queue (e.g., for request prioritization). If this is the case, then several of the above statistics (from “IOdriver Response time stats” to “IOdriver Physical access time stats”) are reported for each of the sub-queues in addition to the above aggregate values. Also, depending upon which sub-queues are employed, up to four additional statistics may be reported:

**IOdriver Requests switched to timeout queue:** the number of requests that were switched to the higher priority *timeout* queue (as described in section 3.4.7), because they were queued in the *base* queue for longer than their specified timeout time.

**IOdriver Timed out requests:** the number of requests that did not complete within their timeout value. Such requests are only switched to the *timeout* queue if they have not yet been initiated.

**IOdriver Half-way timed out requests:** the number of requests that did not complete within half of their timeout value. One of the supported scheduler options gives such requests an intermediate priority level for the remainder of their timeout period.

**IOdriver Requests switched to priority queue:** the number of requests that were switched to the high-priority *timeout* queue because of information delivered from higher-level system components (e.g., the process scheduler) after the request was queued. One reason for such a switch might be that a process must wait for the request to complete [5]; if a high-priority process is waiting on the request completion, the request’s priority may be increased at the I/O driver.

If there is more than one device (or more than one driver), then separate per-driver-per-device statistics are reported. The statistics reported are the same as those described above (as aggregate “IOdriver” statistics). The per-driver-per-device statistics are not reported if “Print driver per-device stats?” is set to false (0).

## 5.2.6 Disk Statistics

The first set of disk statistics is generated by the request queue module. The specific statistics reported are the same as the “IOdriver” statistics described in Section 5.2.5, except that they apply to each disk’s individual request queue(s) (and are denoted accordingly). The “Print ... stats?” parameters for the queue statistics are the same as for the corresponding driver parameters with the word, “driver,” replaced by “disk.”

The next three statistics are not reported if “Print device seek stats?” is set to false (0).

**Disk Seeks of zero distance:** the number of requests resulting in media accesses that require no “seek” (i.e., movement of the disk’s read/write head from one cylinder to another), followed by the fraction of all requests requiring no seek.

**Disk Seek distance stats:** aggregate statistics for the seek distances observed for requests requiring media access.

**Disk Seek time stats:** aggregate statistics for the seek times observed for requests requiring media access.

The next three statistics are not reported if “Print device latency stats?” is set to false (0).

**Disk Full rotation time:** the amount of time required for the disk platters to complete a full revolution. This statistic is only reported for single-disk configurations or in sets of per-disk statistics (see below).

**Disk Zero rotate latency:** the number of media accesses that incur no rotational latency, followed by the fraction of all media accesses incurring no rotational latency.

**Disk Rotational latency stats:** aggregate statistics for the rotational latencies for requests requiring media access. The next statistic is not reported if “Print device xfer stats?” is set to false (0).

**Disk Transfer time stats:** aggregate statistics for the media transfer times for requests requiring media access. The next two statistics are not reported if “Print device acctime stats?” is set to false (0).

**Disk Positioning time stats:** aggregate statistics for positioning times (seek time plus rotational latency) for requests requiring media access.

**Disk Access time stats:** aggregate statistics for media access times for requests requiring media access. The next two statistics are not reported if “Print device interfere stats?” is set to false (0).

**Disk Sequential interference:** the number of requests marked (by a flag) as logically sequential that were not temporally and physically sequential due to interference with other request streams or data mapping algorithms (e.g., striping).

**Disk Local interference:** the number of requests marked (by a flag) as logically “local” that were not temporally or physically local due to interference with other request streams or data mapping algorithms (e.g., striping).

The next seventeen statistics are not reported if “Print device buffer stats?” is set to false (0).

**Disk Number of buffer accesses:** the number of requests that check the disk’s on-board cache for specific contents.

**Disk Buffer hit ratio:** the number of requests that check the disk’s on-board cache and find some “usable” data, followed by the fraction of all requests that check the disk’s on-board cache and find some “usable” data. For example, a read request whose first sector of requested data is in the cache (or is currently being read into the cache) would fall into this category. Also, a write request may fall into this category if the on-board controller allows its data to be appended or prepended to an existing quantity of “dirty” data. In the latter case, the existing dirty data is “usable” because the new request may be combined with it (i.e., is logically sequential to it).

**Disk Buffer miss ratio:** the number of requests that check the disk’s on-board cache and do not find any “usable” data, followed by the fraction of all requests that check the disk’s on-board cache and do not find some “usable” data. For example, a read request whose first sector of requested data is not in the cache (and is not currently being read into the cache) would certainly fall into this category. Also, a write request would fall into this category if the request’s data cannot be combined with any existing “dirty” data in the cache.

**Disk Buffer read hit ratio:** the number of read requests that check the disk’s on-board cache and find all of the requested data already present, followed by the fraction of all read requests and the fraction of all requests that fall into this category.

**Disk Buffer prepend hit ratio:** the number of all write requests that check the disk’s on-board cache and are combined with existing write requests where the new request’s data are logically prepended to the existing “dirty” data, followed by the fraction of all write requests that fall into this category.

**Disk Buffer append hit ratio:** the number of all write requests that check the disk’s on-board cache and are combined with existing write requests where the new request’s data are logically appended to the existing “dirty” data, followed by the fraction of all write requests that fall into this category.

**Disk Write combinations:** the number of all write requests that check the disk’s on-board cache and are combined with existing write requests (either logically prepended or appended), followed by the fraction of all write requests that are combined with existing write requests.

**Disk Ongoing read-ahead hit ratio:** the number of all read requests that check the disk’s on-board cache and find an initial portion of the requested data already present and additional data being actively prefetched into the same cache segment, followed by the fraction of all read requests and the fraction of all requests that fall into this category.

**Disk Average read-ahead hit size:** the average amount of requested data found in the cache for read requests that check the disk’s on-board cache and find an initial portion of the requested data already present and additional data being actively prefetched into the same cache segment.

**Disk Average remaining read-ahead:** the average amount of requested data remaining to be fetched into the cache for read requests that check the disk’s on-board cache and find an initial portion of the requested data already present and additional data being actively prefetched into the same cache segment.

**Disk Partial read hit ratio:** the number of read requests that check the disk's on-board cache and find an initial portion of the requested data already present (with no ongoing prefetch), followed by the fraction of all read requests and the fraction of all requests that fall into this category.

**Disk Average partial hit size:** the average amount of requested data found in the cache for read requests that check the disk's on-board cache and find an initial portion of the requested data (with no ongoing prefetch).

**Disk Average remaining partial:** the average amount of requested data remaining to be fetched into the cache for read requests that check the disk's on-board cache and find an initial portion of the requested data (with no ongoing prefetch).

**Disk Total disk bus wait time:** the total amount of time spent waiting for access to a bus (i.e., arbitration delay).

**Disk Number of disk bus waits:** the total number of times a delay occurred when attempting to access the bus (i.e., the bus was "owned" by another entity when access was requested).

Per-disk statistics are reported for multi-disk configurations. Statistics for specific disks can be enabled or disabled by setting the corresponding "Print stats for disk" configuration parameter (see Section 3.4.5) true (1) or false (0).

### 5.2.7 Controller Statistics

No statistics are reported for the two simple controller models. The following statistics are reported only for "CTLR\_SMART" controllers, which include a cache and are capable of queueing/scheduling requests for one or more attached storage devices. All of the cache statistics are reported for individual controllers only (i.e., no aggregates across controllers are reported). The controller cache statistics are not reported if "Print controller cache stats?" is set to false (0).

**Controller #X cache requests:** the number of requests serviced by the cache of controller X.

**Controller #X cache read requests:** the number of read requests serviced by the cache of controller X, followed by the fraction of serviced requests that are reads.

**Controller #X cache atoms read:** the number of cache atoms accessed by read requests to controller X, followed by the fraction of cache atom accesses that are reads. A "cache atom" is the minimal unit of cache access. In the current version of DiskSim, the cache atom size is always equal to the sector size of the underlying storage devices.

**Controller #X cache read misses:** the number of cache read requests to controller X for which no useful data are found in the cache, followed by the fraction of all requests that are cache read misses and the fraction of all read requests that are misses.

**Controller #X cache read full hits:** the number of cache read requests to controller X for which all necessary data are found in the cache, followed by the fraction of all requests that are cache read full hits and the fraction of all read requests that are full hits.

**Controller #X cache fills (read):** the number of cache fill accesses issued to the underlying storage devices by controller X, followed by the fraction of requests that require a cache fill and the fraction of read requests that require a cache fill.

**Controller #X cache atom fills (read):** the number of atoms read by cache fill accesses issued to the underlying storage devices by controller X, followed by the fraction of cache atom accesses that require a cache fill and the fraction of cache atom read accesses that require a cache fill.

**Controller #X cache write requests:** the number of write requests serviced by the cache of controller X, followed by the fraction of requests that are writes.

**Controller #X cache atoms written:** the number of cache atoms written by write requests to controller X, followed by the fraction of cache atom accesses that are writes.

**Controller #X cache write misses:** the number of cache write requests to controller X that do not overlap at all with data found in the cache, followed by the fraction of all requests that are cache write misses and the fraction of write requests that are misses.

**Controller #X cache write hits (clean):** the number of cache write requests to controller X that overlap only with clean data found in the cache, followed by the fraction of all requests that are clean cache write hits and the fraction of write requests that are clean hits.

**Controller #X cache write hits (dirty):** the number of cache write requests to controller X that overlap with some amount of dirty data found in the cache, followed by the fraction of all requests that are dirty cache write hits and the fraction of write requests that are dirty hits.

**Controller #X cache fills (write):** the number of cache fill accesses (i.e., installation reads [13]) that were required to complete write requests to controller X, followed by the fraction of all requests that require an installation read and

the fraction of all write requests that require an installation read.

**Controller #X cache atom fills (write):** the number of atoms read into the cache in order to complete write requests to controller X, followed by the fraction of all cache atom accesses requiring an installation read and the fraction of all cache atom writes requiring an installation read.

**Controller #X cache destages (write):** the number of destage accesses (i.e., write-backs) initiated by controller X, followed by the fraction of all requests that (eventually) generated a destage access and the fraction of all write requests that generated a destage access.

**Controller #X cache atom destages (write):** the number of atoms written back from the cache of controller X to the storage devices, followed by the fraction of all atom accesses generating (eventually) a destage access and the fraction of all atom write accesses generating a destage access.

**Controller #X cache end dirty atoms:** the number of dirty atoms left in the cache of controller X at the end of the simulation, followed by the fraction of all cache atom accesses that remain dirty at the end of the simulation.

In addition to the per-controller cache statistics, a set of per-controller aggregate queue statistics are generated by the request queue module. That is, queue statistics are reported for each individual controller across all storage devices attached to that controller. The specific statistics reported are the same as the “IOdriver ...” statistics described in Section 5.2.5, except that they apply to each controller’s back-end, per-device request queues (and are denoted accordingly). The “Print ... stats?” parameters for the queue statistics are the same as for the corresponding driver parameters with the word, “driver,” replaced by “controller.”

If there are multiple devices attached to a controller, then the corresponding per-device queue statistics are also reported for each device (i.e., in addition to the aggregate statistics described above). The per-device statistics will not be reported if “Print controller per-device stats?” is set to false (0).

**Total controller bus wait time:** the total amount of time spent by all controllers waiting for access to a bus (i.e., arbitration delay).

### 5.2.8 Bus Statistics

No aggregate statistics (across sets of buses) are reported.

**Bus #X Total utilization time:** the amount of time (in milliseconds) that the bus was not idle during the simulation run. Utilization as a fraction of total simulation time is also reported on this line.

The following set of statistics are not reported if “Print bus idle stats?” is set to false (0).

**Bus #X Idle period length stats:** aggregate statistics for the lengths of idle periods (i.e., periods during which the bus was unused) observed for bus X.

The remaining statistics are not reported if “Print bus arbwait stats?” is set to false (0).

**Bus #X Number of arbitrations:** the number of arbitration decisions made for bus X, including those that involved only a single requester.

**Bus #X Arbitration wait time stats:** aggregate statistics for bus X acquisition delays experienced by attached components. Such delays include both the bus arbitration overhead and any wait time experienced while other components finish their bus transfers.



Parameter	HP C2247A	Seagate ST41601N	DEC RZ26	HP C2490A	HP C3323A
Formatted Capacity	1.05 GB	1.37 GB	1.03 GB	2.13 GB	1.05 GB
RPM	5400	5400	5400	6400	5400
Diameter	3 1/2"	5 1/4"	3 1/2"	3 1/2"	3 1/2"
Height	1.63"	3.25"	1.63"	1.63"	1.00"
Data Surfaces	13	17	14	18	7
Cylinders	2051	2098	2570	2582	2910
Zones	8	14	1	11	8
Sectors/Track	56-96	61-85	58	68-108	72-120

Table 2: Basic disk drive parameters.

## 6 Validation

The disk module of the storage subsystem simulator has been validated by exercising five disk drives, representing three different disk manufacturers (see Table 2), and capturing traces of the resulting I/O activity. Using the observed inter-request delays, each traced request stream was also run through the simulator, which was configured to emulate the corresponding real subsystem. For each disk, this process was repeated for several synthetic workloads with varying read/write ratios, arrival rates, request sizes and degrees of sequentiality and locality. The measured and simulated response time averages match to within 0.8% for all validation runs. (The bus, controller and device driver modules have also been validated as part of a more comprehensive, system-level simulation environment [1].)

Greater insight into the validity of a storage subsystem model can be gained by comparing measured and simulated response time distributions [16]. Figures 3 and 4 show distributions of measured and simulated response times for a sample validation workload of 10,000 requests. Ruemmler and Wilkes[16] define the root mean square horizontal distance between the two distribution curves as a **demerit figure** for disk model calibration. The demerit figure for each of the curves is given in the corresponding caption. The worst-case demerit figure observed over all validation runs was only 2.0% of the corresponding average response time. To our knowledge, no previous disk drive simulator has achieved this level of accuracy.

To accurately mimic the performance behavior of a disk drive, the parameter values used to configure the simulator must accurately reflect the behavior of the actual device. The extremely close match shown in Figure 3 was realized by measuring parameter values directly with a logic analyzer attached to the SCSI bus. The configuration values for the other four disks were obtained with an automatic (software) extraction tool (described in [24]). While still accurate and much less time-consuming, these values are not quite as precise as those obtained with the logic analyzer. We have further observed that using the limited information generally provided in disk drive specifications yields much larger discrepancies between simulated and observed performance.

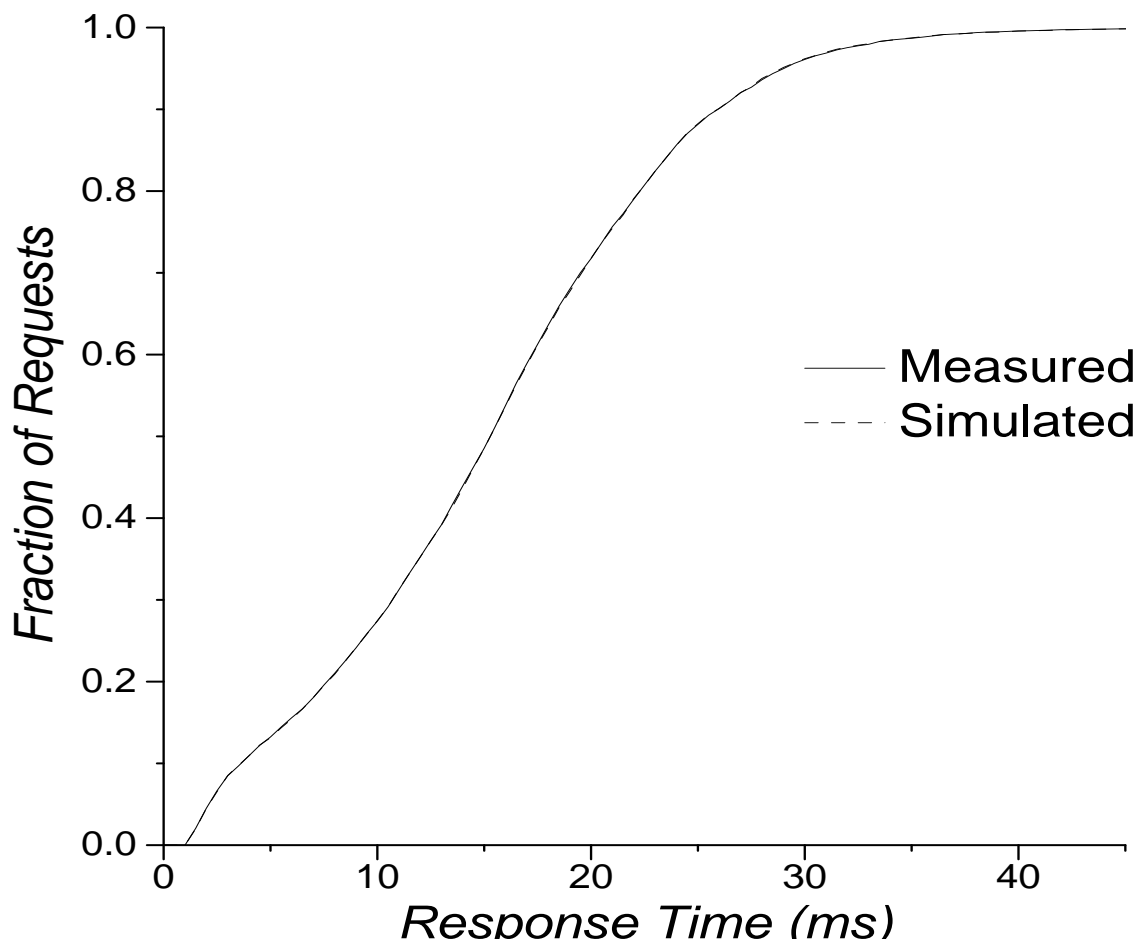
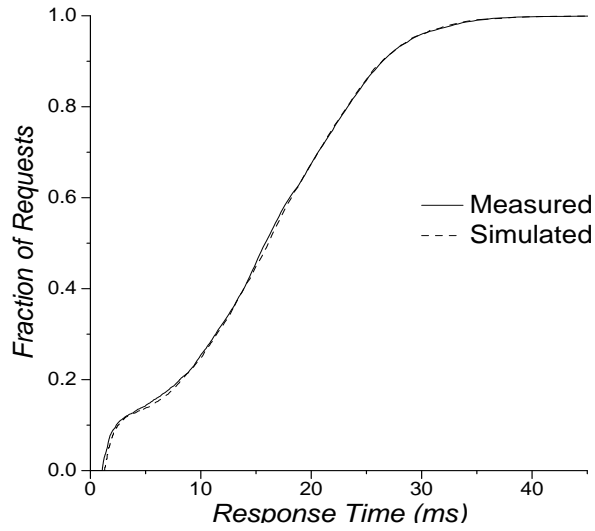
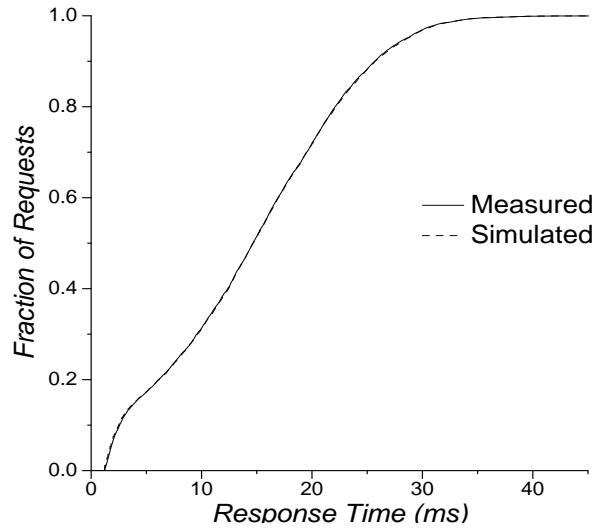


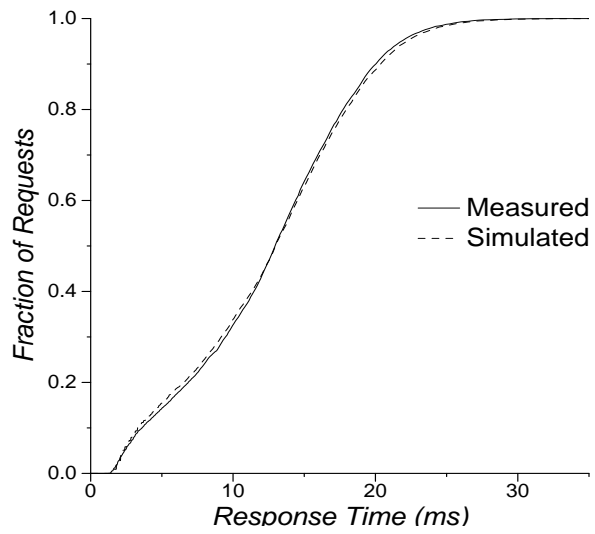
Figure 3: Measured and Simulated Response Time Distributions for an HP C2247A Disk Drive. The demerit figure for this validation run is 0.07 ms, or 0.5% of the corresponding mean response time. Characteristics of the HP C2247A can be found in table 2 and in [8, 22]. The validation workload parameters are 50% reads, 30% sequential, 30% local [normal with 10000 sector variance], 8KB mean request size [exponential], and interarrival time [uniform 0–22 ms].



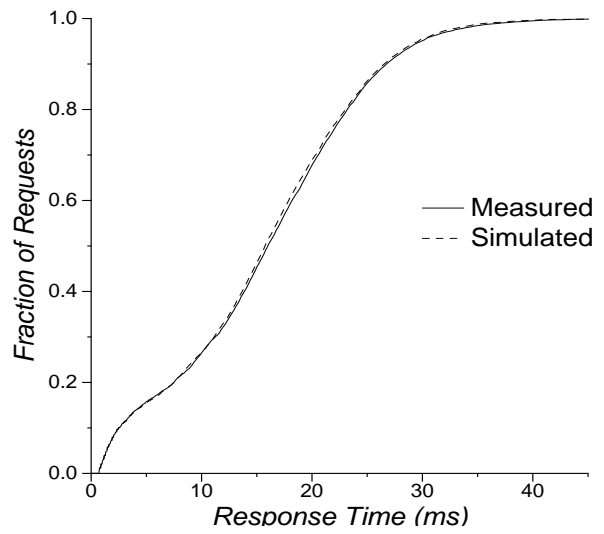
(a) DEC RZ26



(b) Seagate Elite ST41601N



(c) HP C2490A



(d) HP C3323A

Figure 4: Measured and Simulated Response Time Distributions for 4 Disk Drives. The demerit figures for these validation runs are 0.19 ms, 0.075 ms, 0.26 ms and 0.32 ms, respectively (or 1.2%, 0.5%, 2.0% and 1.9% of the corresponding mean response times). Characteristics of these drives can be found in table 2 and in [19, 20, 9, 7, 23]. The validation workload parameters are 50% reads, 30% sequential, 30% local [normal with 10000 sector variance], 8KB mean request size [exponential], and interarrival time [uniform 0–22 ms].

## A Copyright notices for DiskSim

### A.1 Version 3.0 Copyright Addendum

DiskSim Storage Subsystem Simulation Environment (Version 3.0)

Revision Authors: John Bucy, Greg Ganger

Contributors: John Griffin, Jiri Schindler, Steve Schlosser

Copyright (c) of Carnegie Mellon University, 2001, 2002, 2003.

This software is being provided by the copyright holders under the following license. By obtaining, using and/or copying this software, you agree that you have read, understood, and will comply with the following terms and conditions:

Permission to reproduce, use, and prepare derivative works of this software is granted provided the copyright and “No Warranty” statements are included with all reproductions and derivative works and associated documentation. This software may also be redistributed without charge provided that the copyright and “No Warranty” statements are included in all redistributions.

NO WARRANTY. THIS SOFTWARE IS FURNISHED ON AN “AS IS” BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED AS TO THE MATTER INCLUDING, BUT NOT LIMITED TO: WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY OF RESULTS OR RESULTS OBTAINED FROM USE OF THIS SOFTWARE. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT. COPYRIGHT HOLDERS WILL BEAR NO LIABILITY FOR ANY USE OF THIS SOFTWARE OR DOCUMENTATION.

### A.2 Version 2.0 Copyright Addendum

DiskSim Storage Subsystem Simulation Environment (Version 2.0)

Revision Authors: Greg Ganger

Contributors: Ross Cohen, John Griffin, Steve Schlosser

Copyright (c) of Carnegie Mellon University, 1999.

Permission to reproduce, use, and prepare derivative works of this software for internal use is granted provided the copyright and “No Warranty” statements are included with all reproductions and derivative works. This software may also be redistributed without charge provided that the copyright and “No Warranty” statements are included in all redistributions.

NO WARRANTY. THIS SOFTWARE IS FURNISHED ON AN “AS IS” BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED AS TO THE MATTER INCLUDING, BUT NOT LIMITED TO: WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY OF RESULTS OR RESULTS OBTAINED FROM USE OF THIS SOFTWARE. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT.

### A.3 Original (Version 1.0) Copyright Statement

DiskSim Storage Subsystem Simulation Environment

Authors: Greg Ganger, Bruce Worthington, Yale Patt

Copyright (C) 1993, 1995, 1997 The Regents of the University of Michigan

This software is being provided by the copyright holders under the following license. By obtaining, using and/or copying this software, you agree that you have read, understood, and will comply with the following terms and conditions:

Permission to use, copy, modify, distribute, and sell this software and its documentation for any purpose and without fee or royalty is hereby granted, provided that the full text of this NOTICE appears on ALL copies of the

software and documentation or portions thereof, including modifications, that you make.

THIS SOFTWARE IS PROVIDED "AS IS," AND COPYRIGHT HOLDERS MAKE NO REPRESENTATIONS OR WARRANTIES, EXPRESS OR IMPLIED. BY WAY OF EXAMPLE, BUT NOT LIMITATION, COPYRIGHT HOLDERS MAKE NO REPRESENTATIONS OR WARRANTIES OF MERCHANTABILITY OR FITNESS FOR ANY PARTICULAR PURPOSE OR THAT THE USE OF THE SOFTWARE OR DOCUMENTATION WILL NOT INFRINGE ANY THIRD PARTY PATENTS, COPYRIGHTS, TRADEMARKS OR OTHER RIGHTS. COPYRIGHT HOLDERS WILL BEAR NO LIABILITY FOR ANY USE OF THIS SOFTWARE OR DOCUMENTATION.

This software is provided AS IS, WITHOUT REPRESENTATION FROM THE UNIVERSITY OF MICHIGAN AS TO ITS FITNESS FOR ANY PURPOSE, AND WITHOUT WARRANTY BY THE UNIVERSITY OF MICHIGAN OF ANY KIND, EITHER EXPRESSED OR IMPLIED, INCLUDING WITHOUT LIMITATION THE IMPLIED MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE. THE REGENTS OF THE UNIVERSITY OF MICHIGAN SHALL NOT BE LIABLE FOR ANY DAMAGES, INCLUDING SPECIAL, INDIRECT, INCIDENTAL, OR CONSEQUENTIAL DAMAGES, WITH RESPECT TO ANY CLAIM ARISING OUT OF OR IN CONNECTION WITH THE USE OF OR IN CONNECTION WITH THE USE OF THE SOFTWARE, EVEN IF IT HAS BEEN OR IS HEREAFTER ADVISED OF THE POSSIBILITY OF SUCH DAMAGES

The names and trademarks of copyright holders or authors may NOT be used in advertising or publicity pertaining to the software without specific, written prior permission. Title to copyright in this software and any associated documentation will at all times remain with copyright holders.

## B Diskmodel

### B.1 Introduction

Diskmodel is a library implementing mechanical and layout models of modern magnetic disk drives. Diskmodel models two major aspects of disk operation. The layout module models logical-to-physical mapping of blocks, defect management and also computes angular offsets of blocks. The mechanical model handles seek times, rotational latency and various other aspects of disk mechanics.

The implementations of these modules in the current version of Diskmodel are derived from DiskSim 2.0 [4]. DiskSim 3.0 uses Diskmodel natively. Diskmodel has also been used in a device driver implementation of a shortest positioning time first disk request scheduler.

### B.2 Types and Units

All math in diskmodel is performed using integer arithmetic. Angles identified as points on a circle divided into discrete units. Time is represented as multiples of some very small time base. Diskmodel exports the types `dm_time_t` and `dm_angle_t` to represent these quantities. Diskmodel exports functions `dm_time_itod`, `dm_time_dtoi` (likewise for angles) for converting between doubles and the native format. The time function converts to and from milliseconds; the angle function converts to and from a fraction of a circle. `dm_time_t` and `dm_angle_t` should be regarded as opaque and may change over time. Diskmodel is sector-size agnostic in that it assumes that sectors are some fixed size but does not make any assumption about what that size is.

#### B.2.1 Three Zero Angles

When considering the angular offset of a sector on a track, there are at least three plausible candidates for a “zero” angle. The first is “absolute” zero which is the same on every track on the disk. For various reasons, this zero may not coincide with a sector boundary on a track. This motivates the second 0 which we will refer to as  $0_t$  (t for “track”) which is the angular offset of the first sector boundary past 0 on a track. Because of skews and defects, the lowest lbn on the track may not lie at  $0_t$ . We call the angle of the lowest sector on the track  $0_l$  (l for “logical” or “lbn”).

#### B.2.2 Two Zero Sectors

Similarly, when numbering the sectors on a track, it is reasonable to call either the sector at  $0_t$  or the one at  $0_l$  “sector 0.”  $0_t$  corresponds to directly to the physical location of sectors on a track whereas  $0_l$  corresponds to logical layout. Diskmodel works in both systems and the following function descriptions identify which numbering a given function uses.

#### B.2.3 Example

Consider a disk with 100 sectors per track, 2 heads, a head switch skew of 10 sectors and a cylinder switch skew of 20 sectors.  $(x, y, z)$  denotes cylinder  $x$ , head  $y$  and sector  $z$ .

LBN	$0_l$ PBN	$0_t$ PBN
0	(0,0,0)	(0,0,0)
	⋮	
99	(0,0,99)	(0,0,99)
100	(0,1,0)	(0,1,10)
101	(0,1,1)	(0,1,11)
	⋮	
189	(0,1,89)	(0,1,99)
190	(0,1,90)	(0,1,0)
191	(0,1,91)	(0,1,1)
199	(0,1,99)	(0,1,9)

Note that a sector is 3.6 degrees wide.

Cylinder	Head	$\theta_l$ angle
0	0	0 degrees
0	1	36 degrees
1	0	72 degrees
1	1	108 degrees
2	0	180 degrees

## B.3 API

This section describes the data structures and functions that comprise the Diskmodel API.

The `dm_disk_if` struct is the “top-level” handle for a disk in diskmodel. It contains a few disk-wide parameters – number of heads/surfaces, cylinders and number of logical blocks exported by device – along with pointers to the mechanics and layout interfaces.

### B.3.1 Disk-wide Parameters

The top-level of a disk model is the `dm_disk_if` struct:

```
struct dm_disk_if {
    int dm_cyls;           // number of cylinders
    int dm_surfaces;      // number of media surfaces used for data
    int dm_sectors;       // LBNs or total physical sectors (??)

    struct dm_layout_if  *layout;
    struct dm_mech_if    *mech;
};
```

All fields of diskmodel API structures are read-only; the behavior of diskmodel after any of them is modified is undefined. `layout` and `mech` are pointers to the layout and mechanical module interfaces, respectively. Each is a structure containing a number of pointers to functions which constitute the actual implementation. In the following presentation, we write the functions as declarations rather than as types of function pointers for readability. Many of the methods take one or more result parameters; i.e. pointers whose addresses will be filled in with some result. Unless otherwise specified, passing NULL for result parameters is allowed and the result will not be filled in.

### B.3.2 Layout

The layout interface uses the following auxiliary type:

`dm_ptol_result_t` appears in situations where a client code provides a pbn which may not exist on disk as-described e.g. due to defects. It contains the following values:

```
DM_SLIPPED
DM_REMAPPED
DM_OK
DM_NX
```

DM\_SLIPPED indicates that the pbn is a slipped defect. DM\_REMAPPED indicates that the pbn is a remapped defect. DM\_OK indicates that the pbn exists on disk as-is. DM\_NX indicates that there is no sector on the device corresponding to the given pbn. When interpreted as integers, these values are all less than zero so they can be unambiguously intermixed with nonnegative integers e.g. lbns.

The layout module exports the following methods:

```
dm_ptol_result_t dm_translate_ltop(struct dm_disk_if *,
                                int lbn,
                                dm_layout_maptype,
                                struct dm_pbn *result,
                                int *remapsector);
```

Translate a logical block number (lbn) to a physical block number (pbn). `remapsector` is a result parameter which will be set to a non-zero value if the lbn was remapped.

The sector number in the result is relative to the  $0_t$  zero sector.

```
dm_ptol_result_t dm_translate_ltop_0t(struct dm_disk_if *,
                                     int lbn,
                                     dm_layout_maptype,
                                     struct dm_pbn *result,
                                     int *remapsector);
```

Same as `dm_translate_ltop` except that the sector in result is relative to the  $0_t$  sector.

```
dm_ptol_result_t dm_translate_ptol(struct dm_disk_if *,
                                   struct dm_pbn *p,
                                   int *remapsector);
```

Translate a pbn to an lbn. `remapsector` is a result parameter which will be set to a non-zero value if the pbn is defective and remapped.

The sector number in the operand is relative to the  $0_t$  zero sector.

```
dm_ptol_result_t dm_translate_ptol_0t(struct dm_disk_if *,
                                       struct dm_pbn *p,
                                       int *remapsector);
```

Same as `dm_translate_ptol` except that the sector in the result is relative to the  $0_t$  sector.

```
int dm_get_sectors_lbn(struct dm_disk_if *d,
                      int lbn);
```

Returns the number of sectors on the track containing the given lbn.

```
int dm_get_sectors_pbn(struct dm_disk_if *d,
                      struct dm_pbn *);
```

Returns the number of sectors on the track containing the given pbn.

```
void dm_get_track_boundaries(struct dm_disk_if *d,
                             struct dm_pbn *,
                             int *first_lbn,
                             int *last_lbn,
                             int *remapsector);
```

Computes lbn boundaries for the track containing the given pbn. `first_lbn` is a result parameter which returns the first lbn on the track containing the given pbn; similarly, `last_lbn` returns the last lbn on the given track. `remapsector` returns a non-zero value if the first or last block on the track are remapped.

```
dm_ptol_result_t dm_seek_distance(struct dm_disk_if *,
                                  int start_lbn,
                                  int dest_lbn);
```



Computes the seek distance in cylinders that would be incurred for given request. Returns a `dm_ptol_result_t` since one or both of the LBNs may be slipped or remapped.

```
dm_angle_t dm_pbn_skew(struct dm_disk_if *,
                      struct dm_pbn *);
```

This computes the starting offset of a pbn relative to 0. The operand is a pbn relative to  $0_l$ ; the result is an angle relative to 0. This accounts for all skews, slips, etc.

```
dm_angle_t dm_get_track_zerol(struct dm_disk_if *,
                              struct dm_mech_state *);
```

The return value is  $0_l$  for the track identified by the second argument. This is equivalent to calling `dm_pbn_skew` for sector 0 on the same track.

```
dm_ptol_result_t dm_convert_atop(struct dm_disk_if *,
                                 struct dm_mech_state *,
                                 struct dm_pbn *);
```

Finds the pbn of the sector whose leading edge is less than or equal to the given angle. Returns a `ptol_result_t` since the provided angle could be in slipped space, etc. Both the angle in the second operand and the sector number in the result pbn are relative to  $0_l$ .

```
dm_angle_t dm_get_sector_width(struct dm_disk_if *,
                               struct dm_pbn *track,
                               int num);
```

Returns the angular width of an extent of `num` sectors on the given track. Returns 0 if `num` is greater than the number of sectors on the track.

```
dm_angle_t dm_lbn_offset(struct dm_disk_if *, int lbn1, int lbn2);
```

Computes the angular distance/offset between two logical blocks.

```
int dm_marshaled_len(struct dm_disk_if *);
```

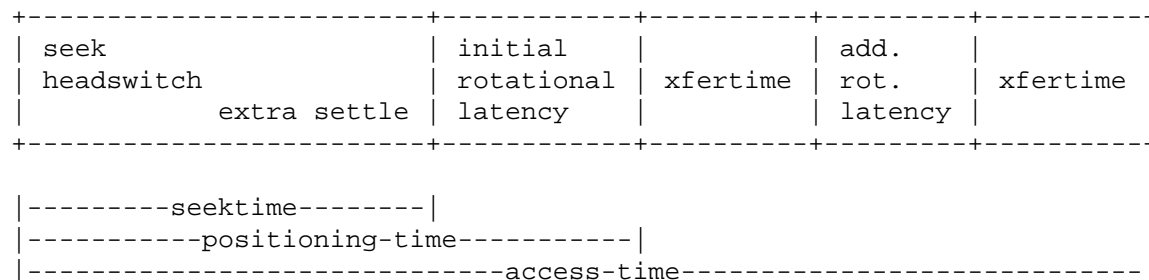
Returns the size of the structure in bytes when marshalled.

```
void *dm_marshall(struct dm_disk_if *, char *);
```

Marshall this layout struct into the provided buffer. The return value is the first address in the buffer not written.

### B.3.3 Mechanics

The following diagram shows the breakdown of a zero-latency access in our model, and the corresponding definitions of seek time, positioning time and access time.



```
dm_time_t dm_seek_time(struct dm_disk_if *,
                      struct dm_mech_state *start_track,
                      struct dm_mech_state *end_track,
                      int read);
```

Computes the amount of time to seek from the first track to the second track, possibly including a head switch and additional write settling time. This is only track-to-track so the angles in the parameters are ignored. `read` should be nonzero if the access on the destination track is a read and zero if it is a write; extra write-settle time is included in the result for writes.

```
int dm_access_block(struct dm_disk_if *,
                  struct dm_mech_state *initial,
                  int start,
                  int len,
                  int immed);
```

From the given initial condition and access, it will return the first block on the track to be read. The access is for `len` sectors starting at physical sector `start` on the same track as `initial`. `immed` indicates if this is an “immediate” or “zero-latency” access; if `immed` is zero, the result will always be the same as `start`.

```
dm_time_t dm_latency(struct dm_disk_if *,
                   struct dm_mech_state *initial,
                   int start,
                   int len,
                   int immed,
                   dm_time_t *addtolatency);
```

This computes the rotational latency incurred from accessing up to `len` blocks from the track starting from angle `initial` and sector `start`. This will access to the end of the track but not wrap around; e.g. for a sequential access that starts on the given track and switches to another, after reaching the end of the first. The return value is the initial rotational latency; i.e. how long before the media transfer for the first block to be read starts. `addtolatency` is a result parameter returning additional rotational latency as defined in the figure above. Note that for non-zero-latency accesses, `addtolatency` will always be zero. Also note that for zero latency accesses, the latency is the amount of time before the media transfer begins for the first sector i.e. the same sector that would be returned by `dm_access_block()`.

`dm_pos_time` and `dm_acctime` optionally return broken-down components of the result via the following struct:

```
struct dm_mech_acctimes {
    dm_time_t seektime;
    dm_time_t initial_latency;
    dm_time_t initial_xfer;
    dm_time_t addl_latency;
    dm_time_t addl_xfer;
};
```

For a zero-latency access, the last two fields will always be zero. `dm_pos_time` only fills in the first two fields; `dm_acctime` fills in all 5.

```
dm_time_t dm_pos_time(struct dm_disk_if *,
                    struct dm_mech_state *initial,
                    struct dm_pbn *start,
                    int len,
                    int rw,
                    int immed);
```

Compute the amount of time before the media transfer for `len` sectors starting at `start` begins starting with the disk mechanics in state `initial`. 0 for `rw` indicates a write, any other value indicates a read. A non-zero value for `immed` indicates a “zero-latency” access. Positioning time is the same as seek time (including head-switch time and any extra write-settle time) plus initial rotational latency.

```
dm_time_t dm_acctime(struct dm_disk_if *,
                    struct dm_mech_state *initial_state,
                    struct dm_pbn *start,
                    int len,
                    int rw,
                    int immed,
                    struct dm_mech_state *result_state);
```

Estimate how long it will take to access `len` sectors starting with pbn `start` with the disk initially in state `initial`. 0 for `rw` indicates a write; any other value indicates a read. A non-zero value for `immed` indicates a “zero-latency” access. `result_state` is a result parameter which returns the mechanical state of the disk when the access completes.

Access time consists of positioning time (above), transfer time and any additional rotational latency not included in the positioning time, e.g. in the middle of a zero-latency access transfer.

`dm_acctime` ignores defects so it yields a smaller-than-correct result when computing access times on tracks with defective sectors. This is deliberate as the handling of defects is a high-level controller function which varies widely.

```
dm_time_t dm_rottime(struct dm_disk_if *,
                    dm_angle_t begin,
                    dm_angle_t end);
```

Compute how long it will take the disk to rotate from the angle in the first position to that in the second position.

```
dm_time_t dm_xfertime(struct dm_disk_if *d,
                     struct dm_mech_state *,
                     int len);
```

Computes the amount of time to transfer `len` sectors to or from the track designated by the second argument. This is computed in terms of `dm_get_sector_width()` and `dm_rottime()` in the obvious way.

```
dm_time_t dm_headswitch_time(struct dm_disk_if *,
                             int h1,
                             int h2);
```

Returns the amount of time to swith from using the first head to the second.

```
dm_angle_t dm_rotate(struct dm_disk_if *,
                    dm_time_t *time);
```

Returns the angle of the media after `time` has elapsed assuming the media started at angle 0.

```
dm_time_t dm_period(struct dm_disk_if *);
```

Returns the rotational period of the media.

```
int dm_marshaled_len(struct dm_disk_if *);
```

Returns the marshalled size of the structure.

```
void *dm_marshall(struct dm_disk_if *, char *);
```

Marshalls the structure into the given buffer. The return value is the first address in the buffer not written.

## B.4 Model Configuration

Diskmodel uses libparam to input the following blocks of parameter data:

```
dm_disk
dm_layout_g1
dm_layout_g1_zone
dm_mech_g1
```

### B.4.1 dm\_disk

The outer `dm_disk` block contains the top-level parameters which are used to fill in the `dm_disk_if` structure. The only valid value for “Layout Model” is a `dm_layout_g1` block and for “Mechanical Model,” a `dm_mech_g1` block.

<code>dm_disk</code>	Block count	int	required
This specifies the number of data blocks. This capacity is exported by the disk (e.g., to a disk array controller). It is not used directly during simulation, but is compared to a similar value computed from other disk parameters. A warning is reported if the values differ.			

<code>dm_disk</code>	Number of data surfaces	int	required
This specifies the number of magnetic media surfaces (not platters!) on which data are recorded. Dedicated servo surfaces should not be counted for this parameter.			

<code>dm_disk</code>	Number of cylinders	int	required
This specifies the number of physical cylinders. All cylinders that impact the logical to physical mappings should be included.			

<code>dm_disk</code>	Mechanical Model	block	required
This block defines the disk’s mechanical model. Currently, the only available implementation is <code>dm_mech_g1</code> .			

<code>dm_disk</code>	Layout Model	block	required
This block defines the disk’s layout model. Currently, the only available implementation is <code>dm_layout_g1</code> .			

### B.4.2 dm\_layout\_g1

The `dm_layout_g1` block provides parameters for a first generation (g1) layout model.

<code>dm_layout_g1</code>	LBN-to-PBN mapping scheme	int	required
This specifies the type of LBN-to-PBN mapping used by the disk. 0 indicates that the conventional mapping scheme is used: LBNs advance along the 0th track of the 0th cylinder, then along the 1st track of the 0th cylinder, thru the end of the 0th cylinder, then to the 0th track of the 1st cylinder, and so forth. 1 indicates that the conventional mapping scheme is modified slightly, such that cylinder switches do not involve head switches. Thus, after LBNs are assigned to the last track of the 0th cylinder, they are assigned to the last track of the 1st cylinder, the next-to-last track of the 1st cylinder, thru the 0th track of the 1st cylinder. LBNs are then assigned to the 0th track of the 2nd cylinder, and so on (“first cylinder is normal”). 2 is like 1 except that the serpentine pattern does not reset at the beginning of each zone; rather, even cylinders are always ascending and odd cylinders are always descending.			

dm_layout_g1	Sparing scheme used	int	required
<p>This specifies the type of sparing used by the disk. Later parameters determine where spare space is allocated. 0 indicates that no spare sectors are allocated. 1 indicates that entire tracks of spare sectors are allocated at the “end” of some or all zones (sets of cylinders). 2 indicates that spare sectors are allocated at the “end” of each cylinder. 3 indicates that spare sectors are allocated at the “end” of each track. 4 indicates that spare sectors are allocated at the “end” of each cylinder and that slipped sectors do not utilize these spares (more spares are located at the “end” of the disk). 5 indicates that spare sectors are allocated at the “front” of each cylinder. 6 indicates that spare sectors are allocated at the “front” of each cylinder and that slipped sectors do not utilize these spares (more spares are located at the “end” of the disk). 7 indicates that spare sectors are allocated at the “end” of the disk. 8 indicates that spare sectors are allocated at the “end” of each range of cylinders. 9 indicates that spare sectors are allocated at the “end” of each zone. 10 indicates that spare sectors are allocated at the “end” of each zone and that slipped sectors do not use these spares (more spares are located at the “end” of the disk).</p>			

dm_layout_g1	Rangesize for sparing	int	required
<p>This specifies the range (e.g., of cylinders) over which spares are allocated and maintained. Currently, this value is relevant only for disks that use “sectors per cylinder range” sparing schemes.</p>			

dm_layout_g1	Skew units	string	optional
<p>This sets the units with which units are input: revolutions or sectors. The “disk-wide” value set here may be overridden per-zone. The default unit is sectors.</p>			

dm_layout_g1	Zones	list	required
<p>This is a list of zone block values describing the zones/bands of the disk.</p>			

The Zones parameter is a list of zone blocks each of which contains the following fields:

dm_layout_g1_zone	First cylinder number	int	required
<p>This specifies the first physical cylinder in the zone.</p>			

dm_layout_g1_zone	Last cylinder number	int	required
<p>This specifies the last physical cylinder in the zone.</p>			

dm_layout_g1_zone	Blocks per track	int	required
<p>This specifies the number of sectors (independent of logical-to-physical mappings) on each physical track in the zone.</p>			

dm_layout_g1_zone	Offset of first block	float	required
<p>This specifies the physical offset of the first logical sector in the zone. Physical sector 0 of every track is assumed to begin at the same angle of rotation. This may be in either sectors or revolutions according to the “Skew units” parameter.</p>			

dm_layout_g1_zone	Skew units	string	optional
<p>Default is sectors. This value overrides any set in the surrounding layout block.</p>			

dm_layout_g1_zone	Empty space at zone front	int	required
<p>This specifies the size of the “management area” allocated at the beginning of the zone for internal data structures. This area can not be accessed during normal activity and is not part of the disk’s logical-to-physical mapping.</p>			

dm_layout_g1_zone	Skew for track switch	float	optional
This specifies the number of physical sectors that are skipped when assigning logical block numbers to physical sectors at a track crossing point. Track skew is computed by the manufacturer to optimize sequential access. This may be in either sectors or revolutions according to the “Skew units” parameter.			

dm_layout_g1_zone	Skew for cylinder switch	float	optional
This specifies the number of physical sectors that are skipped when assigning logical block numbers to physical sectors at a cylinder crossing point. Cylinder skew is computed by the manufacturer to optimize sequential access. This may be in either sectors or revolutions according to the “Skew units” parameter.			

dm_layout_g1_zone	Number of spares	int	required
This specifies the number of spare storage locations – sectors or tracks, depending on the sparing scheme chosen – allocated per region of coverage which may be a track, cylinder, or zone, depending on the sparing scheme. For example, if the sparing scheme is 1, indicating that spare tracks are allocated at the end of the zone, the value of this parameter indicates how many spare tracks have been allocated for this zone.			

dm_layout_g1_zone	slips	list	required
This is a list of lbns for previously detected defective media locations – sectors or tracks, depending upon the sparing scheme chosen – that were skipped-over or “slipped” when the logical-to-physical mapping was last created. Each integer in the list indicates the slipped (defective) location.			

dm_layout_g1_zone	defects	list	required
This list describes previously detected defective media locations – sectors or tracks, depending upon the sparing scheme chosen – that have been remapped to alternate physical locations. The elements of the list are interpreted as pairs wherein the first number is the original (defective) location and the second number indicates the replacement location. Note that these locations will both be either a physical sector number or a physical track number, depending on the sparing scheme chosen.			

### B.4.3 dm\_mech\_g1

The dm\_mech\_g1 block provides parameters for a first generation (g1) mechanical model.

dm_mech_g1	Access time type	string	required
This specifies the method for computing mechanical delays. Legal values are constant which indicates a fixed per-request access time (i.e., actual mechanical activity is not modeled), averageRotation which indicates that seek activity should be modeled but rotational latency is assumed to be equal to one half of a rotation (the statistical mean for random disk access) and trackSwitchPlusRotation which indicates that both seek and rotational activity should be modeled.			

dm_mech_g1	Constant access time	float	optional
Provides the constant access time to be used if the access time type is set to constant.			

dm_mech_g1	Seek type	string	required
This specifies the method for computing seek delays. Legal values are the following: <code>linear</code> indicates that the single-cylinder seek time, the average seek time, and the full-strobe seek time parameters should be used to compute the seek time via linear interpolation. <code>curve</code> indicates that the same three parameters should be used with the seek equation described in [12] (see Section B.4.3). <code>constant</code> indicates a fixed per-request seek time. The <code>Constant seek time</code> parameter must be provided. <code>hpl</code> indicates that the six-value HPL seek equation values parameter (see below) should be used with the seek equation described in [16] (see below). <code>hplplus10</code> indicates that the six-value HPL seek equation values parameter (see below) should be used with the seek equation described in [16] for all seeks greater than 10 cylinders in length. For smaller seeks, use the 10-value <code>First ten seek times</code> parameter (see below) as in [22]. <code>extracted</code> indicates that a more complete seek curve (provided in a separate file) should be used, with linear interpolation used to compute the seek time for unspecified distances. If <code>extracted</code> layout is used, the parameter <code>Full seek curve</code> (below) must be provided.			
dm_mech_g1	Average seek time	float	optional
The mean time necessary to perform a random seek			
dm_mech_g1	Constant seek time	float	optional
For the “constant” seek type (above).			
dm_mech_g1	Single cylinder seek time	float	optional
This specifies the time necessary to seek to an adjacent cylinder.			
dm_mech_g1	Full strobe seek time	float	optional
This specifies the full-strobe seek time (i.e., the time to seek from the innermost cylinder to the outermost cylinder).			
dm_mech_g1	Full seek curve	string	optional
The name of the input file containing the seek curve data. The format of this file is described below.			
dm_mech_g1	Add. write settling delay	float	required
This specifies the additional time required to precisely settle the read/write head for writing (after a seek or head switch). As this parameter implies, the seek times computed using the above parameter values are for read access.			
dm_mech_g1	Head switch time	float	required
This specifies the time required for a head switch (i.e., activating a different read/write head in order to access a different media surface).			
dm_mech_g1	Rotation speed (in rpms)	int	required
This specifies the rotation speed of the disk platters in rpms.			
dm_mech_g1	Percent error in rpms	float	required
This specifies the maximum deviation in the rotation speed specified above. During initialization, the rotation speed for each disk is randomly chosen from a uniform distribution of the specified rotation speed $\pm$ the maximum allowed error. This feature may be deprecated and should be avoided.			
dm_mech_g1	First ten seek times	list	optional
This is a list of ten floating-point numbers specifying the seek time for seek distances of 1 through 10 cylinders.			

dm_mech_g1	HPL seek equation values	list	optional
This is a list containing six numbers specifying the variables $V_1$ through $V_6$ of the seek equation described in [16] (see below).			

### Lee's Seek Equation

$$seekTime(x) = \begin{cases} 0 & : \text{if } x = 0 \\ a\sqrt{x-1} + b(x-1) + c & : \text{if } x > 0 \end{cases}, \text{ where}$$

$x$  is the seek distance in cylinders,

$$a = (-10minSeek + 15avgSeek - 5maxSeek)/(3\sqrt{numCyl}),$$

$$b = (7minSeek - 15avgSeek + 8maxSeek)/(3numCyl), \text{ and}$$

$$c = minSeek.$$

### The HPL Seek Equation

Seek distance	Seek time
1 cylinder	$V_6$
$<V_1$ cylinders	$V_2 + V_3 * \sqrt{dist}$
$\geq V_1$ cylinders	$V_4 + V_5 * dist$

, where  $dist$  is the seek distance in cylinders.

If  $V_6 == -1$ , single-cylinder seeks are computed using the second equation.  $V_1$  is specified in cylinders, and  $V_2$  through  $V_6$  are specified in milliseconds.

$V_1$  must be a non-negative integer,  $V_2 \dots V_5$  must be non-negative floats and  $V_6$  must be either a non-negative float or  $-1$ .

### Format of an Extracted Seek Curve

An extracted seek file contains a number of (seek-time,seek-distance) data points. The format of such a file is very simple: the first line is

Seek distances measured: <n>

where <n> is the number of seek distances provided in the curve. This line is followed by <n> lines of the form <distance>, <time> where <distance> is the seek distance measured in cylinders, and <time> is the amount of time the seek took in milliseconds. e.g.

Seek distances measured: 4

1, 1.2

2, 1.5

5, 5

10, 9.2



## References

- [1] Gregory R. Ganger. Generating representative synthetic workloads: an unsolved problem. *International Conference on Management and Performance Evaluation of Computer Systems* (Nashville, TN), pages 1263–1269, 1995.
- [2] Gregory R. Ganger. *System-oriented evaluation of I/O subsystem performance*. PhD thesis, published as CSE-TR-243-95. University of Michigan, Ann Arbor, MI, June 1995.
- [3] Gregory R. Ganger and Yale N. Patt. The process-flow model: examining I/O performance from the system's point of view. *ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems*, pages 86–97, May 1993.
- [4] Gregory R. Ganger, Bruce L. Worthington, and Yale N. Patt. *The DiskSim simulation environment version 2.0 reference manual*, December 1999.
- [5] Gregory Robert Ganger. Improved methodologies for evaluating I/O architectures. Electrical Engineering and Computer Science: Computer Science and Engineering Division, University of Michigan, December 1993.
- [6] John Linwood Griffin, Jiri Schindler, Steven W. Schlosser, John C. Bucy, and Gregory R. Ganger. Timing-accurate storage emulation. *Conference on File and Storage Technologies* (Monterey, CA, 28–30 January 2002), pages 75–88. USENIX Association, 2002.
- [7] Hewlett-Packard Company. *HP C3323A 3.5-inch SCSI-2 Disk Drives, Technical Reference Manual Part Number 5962-6452*, second edition, April 1994.
- [8] Hewlett-Packard Company. *HP C2244/45/46/47 3.5-inch SCSI-2 Disk Drive Technical Reference Manual Part Number 5960-8346*, third edition, September 1992.
- [9] Hewlett-Packard Company. *HP C2490A 3.5-inch SCSI-2 Disk Drives, Technical Reference Manual Part Number 5961-4359*, third edition, September 1993.
- [10] Ramakrishna Karedla, J. Spencer Love, and Bradley G. Wherry. Caching strategies to improve disk performance. *IEEE Computer*, **27**(3):38–46, March 1994.
- [11] Edward K. Lee and Randy H. Katz. An analytic performance model of disk arrays. *ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems* (Santa Clara, CA, 17–21 May 1993). Published as *Performance Evaluation Review*, **21**(1):98–109, June 1993.
- [12] Edward Kihyen Lee. *Performance modeling and analysis of disk arrays*. PhD thesis, published as UCB//CSD-93-770. Department of Electrical Engineering and Computer Science, University of California at Berkeley, 1993.
- [13] James O'Toole and Liuba Shrira. Opportunistic log: efficient installation reads in a reliable storage server. *Symposium on Operating Systems Design and Implementation* (Monterey, CA), pages 39–48. Usenix Association, 14–17 November 1994.
- [14] Mendel Rosenblum, Edouard Bugnion, Stephen Alan Herrod, Emmett Witchel, and Anoop Gupta. The impact of architectural trends on operating system performance. *ACM Symposium on Operating System Principles* (Copper Mountain Resort, CO, 3–6 December 1995). Published as *Operating Systems Review*, **29**(5), 1995.
- [15] Chris Ruemmler and John Wilkes. UNIX disk access patterns. *Winter USENIX Technical Conference* (San Diego, CA, 25–29 January 1993), pages 405–420, 1993.
- [16] Chris Ruemmler and John Wilkes. An introduction to disk drive modeling. *IEEE Computer*, **27**(3):17–28, March 1994.
- [17] M. Satyanarayanan. *Modelling storage systems*. UMI Research Press, 1986.
- [18] Jiri Schindler and Gregory R. Ganger. *Automated disk drive characterization*. Technical report CMU-CS-99-176. Carnegie-Mellon University, Pittsburgh, PA, December 1999.
- [19] Seagate Technology, Inc. *SCSI Interface Specification, Small Computer System Interface (SCSI), Elite Product Family Document Number 64721702*, revision D, March 1992.
- [20] Seagate Technology, Inc. *Seagate Product Specification, ST41600N and ST41601N Elite Disc Drive, SCSI Interface Document Number 64403103*, revision G, October 1992.
- [21] Bruce L. Worthington. *Aggressive centralized and distributed scheduling of disk requests*. PhD thesis, published as CSE-TR-244-95. Department of Computer Science and Engineering, University of Michigan, June 1995.
- [22] Bruce L. Worthington, Gregory R. Ganger, and Yale N. Patt. *Scheduling for modern disk drives and non-random workloads*. CSE-TR-194-94. Department of Computer Science and Engineering, University of Michigan, 1 March 1994.
- [23] Bruce L. Worthington, Gregory R. Ganger, Yale N. Patt, and John Wilkes. *On-line extraction of SCSI disk drive parameters*. CSE-TR-323-96. Department of Electrical Engineering and Computer Science, University of Michigan, December 1996.
- [24] Bruce L. Worthington, Gregory R. Ganger, Yale N. Patt, and John Wilkes. On-line extraction of SCSI disk drive parameters. *ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems* (Ottawa, Canada), pages 146–156, May 1995.