

Oracle9i Application Server

Oracle HTTP Server *powered by Apache* Performance Guide

Release 1.0.2.2 for Sun SPARC Solaris

May 2001

Part No. A86059-02

Oracle9i Application Server

Oracle HTTP Server *powered by Apache* Performance Guide, Release 1.0.2.2 for Windows NT

Part No. A86059-02

Copyright © 2001, Oracle Corporation. All rights reserved.

Contributors: Alice Chan, Gary Hallmark, Bruce Irvin, Alexander Hoefling, Sharon Malek, Carol Orange, Mukul Paithane, Leela Rao, Joan Silverman, Sanjay Singh, Eddy So

The Programs (which include both the software and documentation) contain proprietary information of Oracle Corporation; they are provided under a license agreement containing restrictions on use and disclosure and are also protected by copyright, patent, and other intellectual and industrial property laws. Reverse engineering, disassembly, or decompilation of the Programs is prohibited.

The information contained in this document is subject to change without notice. If you find any problems in the documentation, please report them to us in writing. Oracle Corporation does not warrant that this document is error free. Except as may be expressly permitted in your license agreement for these Programs, no part of these Programs may be reproduced or transmitted in any form or by any means, electronic or mechanical, for any purpose, without the express written permission of Oracle Corporation.

If the Programs are delivered to the U.S. Government or anyone licensing or using the programs on behalf of the U.S. Government, the following notice is applicable:

Restricted Rights Notice Programs delivered subject to the DOD FAR Supplement are "commercial computer software" and use, duplication, and disclosure of the Programs, including documentation, shall be subject to the licensing restrictions set forth in the applicable Oracle license agreement. Otherwise, Programs delivered subject to the Federal Acquisition Regulations are "restricted computer software" and use, duplication, and disclosure of the Programs shall be subject to the restrictions in FAR 52.227-19, Commercial Computer Software - Restricted Rights (June, 1987). Oracle Corporation, 500 Oracle Parkway, Redwood City, CA 94065.

The Programs are not intended for use in any nuclear, aviation, mass transit, medical, or other inherently dangerous applications. It shall be the licensee's responsibility to take all appropriate fail-safe, backup, redundancy, and other measures to ensure the safe use of such applications if the Programs are used for such purposes, and Oracle Corporation disclaims liability for any damages caused by such use of the Programs.

Oracle is a registered trademark, and the Oracle Logo, Oracle9i Application Server, Oracle8i, Oracle9i, Oracle Enterprise Manager, Oracle Internet Directory, and PL/SQL are trademarks or registered trademarks of Oracle Corporation. Other names may be trademarks of their respective owners.

This product includes software developed by the Apache Group for use in the Apache HTTP server project (<http://www.apache.org/>).

This product includes software developed by the OpenSSL project for use in the OpenSSL Toolkit (<http://www.openssl.org/>). This product includes cryptographic software written by Eric Young (ey@cryptsoft.com). This product includes software written by Tim Hudson (tjh@cryptsoft.com).

This product includes software developed by Ralf S. Engelschall (rse@engelschall.com) for use in the mod_ssl project (<http://www.modssl.org/>).

Contents

Send Us Your Comments	vii
Preface.....	ix
Audience	x
Organization.....	x
Related Documentation	xi
Conventions.....	xii
1 Performance Overview	
Performance Terms	1-2
What is Performance Tuning?	1-3
System Throughput.....	1-4
Wait Time.....	1-5
Critical Resources	1-5
Effects of Excessive Demand.....	1-7
Adjustments to Relieve Problems	1-7
Setting Performance Targets.....	1-8
Setting User Expectations.....	1-8
Evaluating Performance	1-8
Performance Methodology.....	1-9
Factors in Improving Performance	1-10
Architecture.....	1-11

2 Monitoring Your Web Server

Monitoring Processor Use	2-2
Using the sar Utility	2-2
Using the mpstat Utility	2-3
Monitoring Network Traffic	2-4
Using the snoop Utility	2-4
Monitoring the Web Server	2-6
Using the mod_status Utility to Monitor the Web Server	2-6
Logging Server Statistics to a File	2-9
Monitoring JServ Processes	2-10

3 Sizing and Configuration

Sizing your Hardware and Resources	3-2
Determining CPU Requirements	3-3
Secure Sockets Layer Impact on CPU Requirements	3-4
Determining Memory Requirements	3-4
Determining Memory Requirements for Operating System Software	3-5
Determining Memory Requirements for the Oracle HTTP Server	3-5
JServ Memory Requirements	3-5
Determining Java Heap Size	3-5
Determining Memory Requirements for Servlets and OracleJSP pages	3-6
Determining the Number of JServ Processes per CPU	3-8

4 Optimizing HTTP Server Performance

Tuning TCP Parameters	4-2
Setting TCP parameters	4-2
Configuring the MaxClients Parameter	4-6
Enabling SSL Session Caching	4-7
Understanding Performance Implications of Logging	4-7
Benefits of the HTTP/1.1 Protocol	4-8
Supporting Persistent Connections	4-8
Differences between Apache Releases with Respect to Performance	4-11

5 Optimizing Apache JServ

Overview of JServ	5-2
Optimizing Servlet Performance	5-3
Loading Servlet Classes	5-3
Reloading Servlet Classes Automatically	5-3
How to Perform Load Balancing.....	5-4
Using Single Thread Model Servlets.....	5-7
What is OracleJSP?	5-8
Tuning OracleJSP Pages for Performance	5-8
Impact of Session Management	5-8
Developer Mode	5-9
Buffering	5-9
OracleJSP Performance Tips	5-10

Index

Send Us Your Comments

**Oracle9i Application Server Oracle HTTP Server *powered by Apache* Performance Guide,
Release 1.0.2.2 for Windows NT**

Part No. A86676-02

Oracle Corporation welcomes your comments and suggestions on the quality and usefulness of this publication. Your input is an important part of the information used for revision.

- Did you find any errors?
- Is the information clearly presented?
- Do you need more information? If so, where?
- Are the examples correct? Do you need more examples?
- What features did you like most about this manual?

If you find any errors or have any other suggestions for improvement, please indicate the chapter, section, and page number (if available). You can send comments to us in the following ways:

- Electronic mail - iasdocs_us@oracle.com
- Fax - (650) 506-7409 Attn: Oracle9i Application Server Documentation Manager
- Postal service:

Oracle Corporation
Oracle9i Application Server Documentation Manager
500 Oracle Parkway, M/S 2op4
Redwood Shores, CA 94065 USA

If you would like a reply, please give your name, address, and telephone number below.

If you have problems with the software, please contact your local Oracle Support Services.

Preface

This guide discusses configuration and performance tuning of the Oracle HTTP Server *powered by Apache*.

There are many sources of information on configuring and tuning web servers, Apache in particular. This guide refers to those sources when expedient, and, where practical, quantifies the performance gains resulting from configuration actions found in those sources. Any recommendations not validated by Oracle in-house testing are cited as such, with attribution to the original source.

All in-house tests detailed in this guide were run on a dedicated 100 Mbps network, in order to achieve repeatable test results. Your results will vary based on network configuration and contention characteristics.

This preface contains these topics:

- [Audience](#)
- [Organization](#)
- [Related Documentation](#)
- [Conventions](#)

Audience

This guide is written for Oracle9i Application Server developers and system administrators who are responsible for configuring and tuning the Oracle HTTP Server *powered by Apache*.

To use this document, you need a working knowledge of web server administration and performance tuning concepts.

Organization

This document contains:

[Chapter 1, "Performance Overview"](#)

Describes performance and tuning concepts and terminology, with a description of the Oracle HTTP Server components in the Oracle9i Application Server architecture.

[Chapter 2, "Monitoring Your Web Server"](#)

Discusses the importance of monitoring to performance tuning, and describes tools and processes for gathering information about the web server and operating system software.

[Chapter 3, "Sizing and Configuration"](#)

Provides guidelines and approaches to sizing and configuration to meet performance goals.

[Chapter 4, "Optimizing HTTP Server Performance"](#)

Discusses tuning parameters to improve HTTP server performance and the effects of caching and logging on performance.

[Chapter 5, "Optimizing Apache JServ"](#)

Discusses performance and load balancing for Apache JServ, and optimizing the performance of OracleJSP pages.

Related Documentation

For more information, see these Oracle resources:

- *Oracle9i Application Server Overview Guide*
- *OracleJavaServer Pages Developer's Guide and Reference*

In North America, printed documentation is available for sale in the Oracle Store at

<http://oraclestore.oracle.com/>

Customers in Europe, the Middle East, and Africa (EMEA) can purchase documentation from

<http://www.oraclebookshop.com/>

Other customers can contact their Oracle representative to purchase printed documentation.

To download free release notes, installation documentation, white papers, or other collateral, please visit the Oracle Technology Network (OTN). You must register online before using OTN; registration is free and can be done at

<http://technet.oracle.com/membership/index.htm>

If you already have a username and password for OTN, then you can go directly to the documentation section of the OTN Web site at

<http://technet.oracle.com/docs/index.htm>

The following sources provide additional information on topics found in this guide:

- *Sun Performance and Tuning: Java and the Internet* by Adrian Cockcroft and Richard Petit, Prentice Hall, 1998.
- *Solaris 2.X: Performance Management, Fine Tuning, & Capacity Planning* by H. Frank Cervone, McGraw Hill, 1998.
- For information on the `mod_status` utility, see
<http://www.oreillynet.com/pub/a/apache/2000/04/21/wrangler.html>
http://www.apache.org/docs/mod/mod_status.html
- For information on the `LogLevel` directive, see
<http://www.apache.org/docs/mod/core.html#loglevel>

- For a discussion on memory usage in Solaris, see the white paper entitled “The Solaris Memory System: Sizing, Tools and Architecture” at <http://www.sun.com/sun-on-net/performance/vmsizing.pdf>
- For information on Apache web server performance, see Dale Gaudet’s *Apache Performance Notes* at <http://www.apache.org/docs/misc/perf-tuning.html>
- For information about performance and the HTTP/1.1 protocol, see <http://www.w3.org/Protocols/HTTP/Performance/Pipeline.html>
- For more information on the FIN_WAIT_2 state, see http://apache.put.poznan.pl/misc/fin_wait_2.html

Conventions

This section describes the conventions used in the text and code examples of the this documentation set. It describes:

- [Conventions in Text](#)
- [Conventions in Code Examples](#)

Conventions in Text

We use various conventions in text to help you more quickly identify special terms. The following table describes those conventions and provides examples of their use.

Convention	Meaning	Example
Bold	Bold typeface indicates terms that are defined in the text or terms that appear in a glossary, or both.	The C datatypes such as ub4 , sword , or OCINumber are valid. When you specify this clause, you create an index-organized table .
<i>Italics</i>	Italic typeface indicates book titles, emphasis, syntax clauses, or placeholders.	<i>Oracle8i Concepts</i> You can specify the <i>parallel_clause</i> . Run <code>Uold_release.SQL</code> where <i>old_release</i> refers to the release you installed prior to upgrading.

Convention	Meaning	Example
lowercase monospace (fixed-width font)	Lowercase monospace typeface indicates executables and sample user-supplied elements. Such elements include computer and database names, net service names, and connect identifiers, as well as user-supplied database objects and structures, column names, packages and classes, user names and roles, program units, and parameter values.	Enter <code>sqlplus</code> to open SQL*Plus. The <code>department_id</code> , <code>department_name</code> , and <code>location_id</code> columns are in the <code>hr.departments</code> table. Set the <code>QUERY_REWRITE_ENABLED</code> initialization parameter to <code>true</code> . Connect as <code>oe</code> user.

Conventions in Code Examples

Code examples illustrate SQL, PL/SQL, SQL*Plus, or other command-line statements. They are displayed in a monospace (fixed-width) font and separated from normal text as shown in this example:

```
SELECT username FROM dba_users WHERE username = 'MIGRATE';
```

The following table describes typographic conventions used in code examples and provides examples of their use.

Convention	Meaning	Example
[]	Brackets enclose one or more optional items. Do not enter the brackets.	<code>DECIMAL (digits [, precision])</code>
{ }	Braces enclose two or more items, one of which is required. Do not enter the braces.	<code>{ENABLE DISABLE}</code>
	A vertical bar represents a choice of two or more options within brackets or braces. Enter one of the options. Do not enter the vertical bar.	<code>{ENABLE DISABLE}</code> <code>[COMPRESS NOCOMPRESS]</code>
...	Horizontal ellipsis points indicate either: <ul style="list-style-type: none"> That we have omitted parts of the code that are not directly related to the example That you can repeat a portion of the code 	<code>CREATE TABLE ... AS subquery;</code> <code>SELECT col1, col2, ... , coln FROM employees;</code>
.	Vertical ellipsis points indicate that we have omitted several lines of code not directly related to the example.	

Convention	Meaning	Example
Other notation	You must enter symbols other than brackets, braces, vertical bars, and ellipsis points as it is shown.	<pre>acctbal NUMBER(11,2); acct CONSTANT NUMBER(4) := 3;</pre>
<i>Italics</i>	Italicized text indicates variables for which you must supply particular values.	<pre>CONNECT SYSTEM/<i>system_password</i></pre>
UPPERCASE	Uppercase typeface indicates elements supplied by the system. We show these terms in uppercase in order to distinguish them from terms you define. Unless terms appear in brackets, enter them in the order and with the spelling shown. However, because these terms are not case sensitive, you can enter them in lowercase.	<pre>SELECT last_name, employee_id FROM employees; SELECT * FROM USER_TABLES; DROP TABLE hr.employees;</pre>
lowercase	Lowercase typeface indicates programmatic elements that you supply. For example, lowercase indicates names of tables, columns, or files.	<pre>SELECT last_name, employee_id FROM employees; sqlplus hr/hr</pre>

Performance Overview

This chapter discusses performance and tuning concepts, and briefly describes Oracle9i Application Server architecture.

This chapter contains the following sections:

- [Performance Terms](#)
- [What is Performance Tuning?](#)
- [Setting Performance Targets](#)
- [Setting User Expectations](#)
- [Evaluating Performance](#)
- [Performance Methodology](#)
- [Architecture](#)

Performance Terms

Following are performance terms used in this book:

concurrency	The ability to handle multiple requests simultaneously. Threads and processes are examples of concurrency mechanisms.
contention	Competition for resources.
hash	A number generated from a string of text with an algorithm. The hash value is substantially smaller than the text itself. Hash numbers are used for security and for faster access to data.
latency	The time that one system component spends waiting for another component in order to complete the entire task. Latency can be defined as wasted time. In networking contexts, latency is defined as the travel time of a packet from source to destination.
response time	The time between the submission of a request and the receipt of the response.
scalability	<p>The ability of a system to provide throughput in proportion to, and limited only by, available hardware resources.</p> <p>A scalable system is one that can handle increasing numbers of requests without adversely affecting response time and throughput.</p>
service time	The time between the receipt of a request and the completion of the response to the request.
think time	The time the user is not engaged in actual use of the processor.
throughput	The number of requests processed per unit of time.
wait time	The time between the submission of the request and initiation of the request.

What is Performance Tuning?

Performance must be built in. You must anticipate performance requirements during application analysis and design, and balance the costs and benefits of optimal performance. This section introduces some fundamental concepts:

- Response Time
- System Throughput
- Wait Time
- Critical Resources
- Effects of Excessive Demand
- Adjustments to Relieve Problems

See Also: [“Setting Performance Targets” on page 1-8](#) for a discussion on performance requirements and determining what parts of the system to tune.

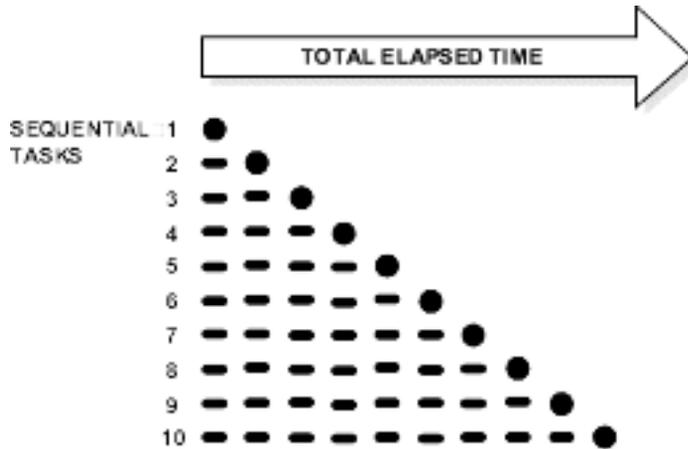
- Response Time

Because **response time** equals **service time** plus **wait time**, you can increase performance in this area by:

- Reducing **wait time**
- Reducing **service time**

[Figure 1-1](#) illustrates ten independent tasks competing for a single resource.

Figure 1–1 Sequential processing of independent tasks



In this example, only task 1 runs without waiting. Task 2 must wait until task 1 has completed; task 3 must wait until tasks 1 and 2 have completed, and so on. (Although the figure shows the independent tasks as the same size, the size of the tasks will vary.)

In parallel processing with multiple resources, more resources are available to the tasks. Each independent task executes immediately using its own resource: no **wait time** is involved.

The Oracle HTTP Server processes requests in this fashion, allocating client requests to available httpd processes. The MaxClients parameter specifies the number of httpd processes simultaneously available to handle client requests. When the number of processes in use reaches the MaxClients value, the server refuses connections until requests are completed and processes are freed.

System Throughput

System **throughput** is the amount of work accomplished in a given amount of time. You can increase **throughput** by:

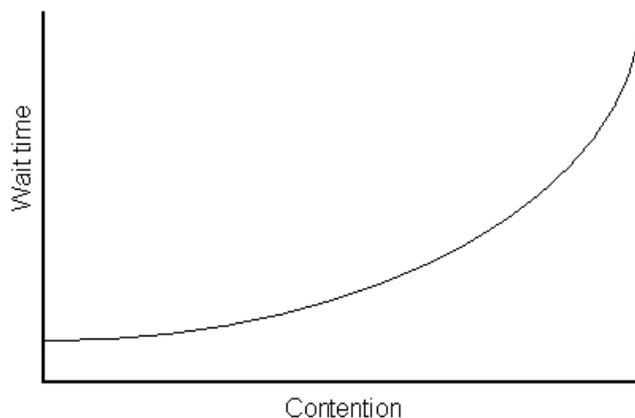
- Reducing **service time**

- Reducing overall **response time** by increasing the amount of scarce resources available. For example, if the system is CPU bound, and you can add more CPUs.

Wait Time

While the **service time** for a task may stay the same, **wait time** will lengthen with increased **contention**. If many users are waiting for a service that takes one second, the tenth user must wait 9 seconds. [Figure 1-2](#) shows the relationship between **wait time** and resource **contention**.

Figure 1-2 *Wait time rising with increased contention for a resource*



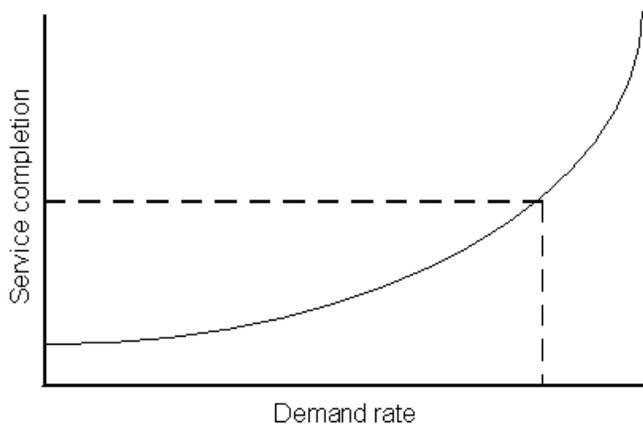
Critical Resources

Resources such as CPU, memory, I/O capacity, and network bandwidth are key to reducing **service time**. Adding resources increases **throughput** and reduces **response time**. Performance depends on these factors:

- How many resources are available?
- How many clients need the resource?
- How long must they wait for the resource?
- How long do they hold the resource?

[Figure 1-3](#) shows that as the number of units requested rises, the time to service completion rises.

Figure 1–3 *Time to service completion vs. demand rate*



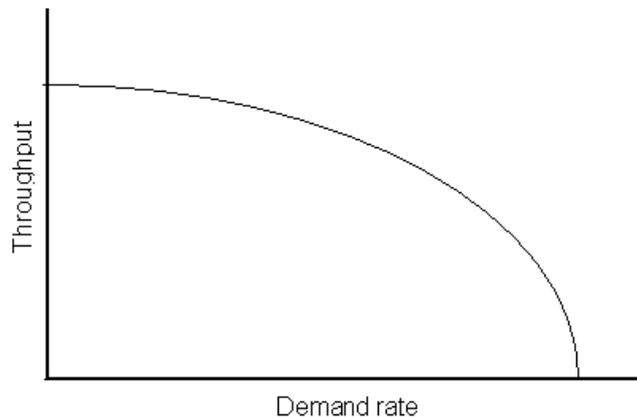
To manage this situation, you have two options:

- Limit demand rate to maintain acceptable **response times**
- Add resources

Effects of Excessive Demand

Excessive demand increases **response time** and reduces **throughput**, as shown in [Figure 1-4](#). If there is any possibility of the demand rate exceeding the achievable **throughput**, then determine which parameters should be adjusted (such as `asMaxClients` in the Oracle HTTP Server and `security.maxConnections` in JServ) and change the configuration accordingly.

Figure 1-4 *Increased Demand/Reduced Throughput*



Adjustments to Relieve Problems

Performance problems can be relieved by making adjustments in the following areas:

unit consumption	Reducing the resource (CPU, memory) consumption of each request can improve performance. This might be achieved by pooling and caching.
functional demand	Rescheduling or redistributing the work will relieve some problems.
capacity	Increasing or reallocating resources (such as CPUs) relieves some problems.

Setting Performance Targets

Whether you are designing or maintaining a system, you should set specific performance goals so that you know how and what to optimize. If you alter parameters without a specific goal in mind, you can waste time tuning your system without significant gain.

An example of a specific performance goal is an order entry **response time** under three seconds. If the application does not meet that goal, identify the cause (for example, I/O **contention**), and take corrective action. During development, test the application to determine if it meets the designed performance goals.

Tuning usually involves a series of trade-offs. Once you have determined the bottlenecks, you may have to modify performance in some other areas to achieve the desired results. For example, if I/O is a problem, you may need to purchase more memory or more disks. If a purchase is not possible, you may have to limit the **concurrency** of the system to achieve the desired performance. However, if you have clearly defined goals for performance, the decision on what to trade for higher performance is simpler because you have identified the most important areas.

Setting User Expectations

Application developers, database administrators, and system administrators must be careful to set appropriate performance expectations for users. When the system carries out a particularly complicated operation, **response time** may be slower than when it is performing a simple operation. Users should be made aware of which operations might take longer.

Evaluating Performance

With clearly defined performance goals, you can readily determine when performance tuning has been successful. Success depends on the functional objectives you have established with the user community, your ability to measure whether or not the criteria are being met, and your ability to take corrective action to overcome any exceptions.

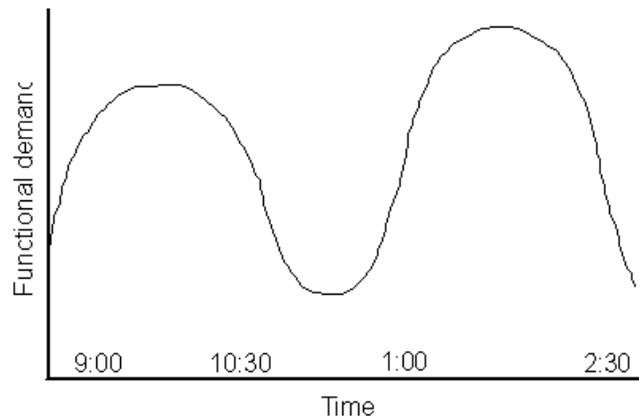
Ongoing performance monitoring enables you to maintain a well tuned system. Keeping a history of the application's performance over time enables you to make useful comparisons. With data about actual resource consumption for a range of loads, you can conduct objective **scalability** studies and from these predict the resource requirements for anticipated load volumes.

Performance Methodology

Achieving optimal effectiveness in your system requires planning, monitoring, and periodic adjustment. The first step in performance tuning is to determine the goals you need to achieve and to design effective usage of available technology into your applications. After implementing your system, it is necessary to periodically monitor and adjust your system. For example, you might want to ensure that 90% of the users experience **response times** no greater than 5 seconds and the maximum **response time** for all users is 20 seconds. Usually, it's not that simple. Your application may include a variety of operations with differing characteristics and acceptable response times. You will need to set measurable goals for each of these.

You will also need to determine variances in the load. For example, users might access the system heavily between 9:00am and 10:00am and then again between 1:00pm and 2:00pm, as shown in [Figure 1-5](#). If your peak load occurs on a regular basis, for example, daily or weekly, the conventional wisdom is to configure and tune systems to meet your peak load requirements. The lucky users who access the application in off-time will experience better **response times** than your peak-time users. If your peak load is infrequent, you may be willing to tolerate higher **response times** at peak loads for the cost savings of smaller hardware configurations.

Figure 1-5 *Adjusting Capacity and Functional Demand*



Factors in Improving Performance

Performance spans several areas:

- **Application design:** Designing applications that efficiently utilize hardware resources and handle increasing numbers of users effectively.
- **Sizing and configuration:** Determining the type of hardware needed to support your performance goals. See [Chapter 3, "Sizing and Configuration"](#).
- **Parameter tuning:** Setting configurable parameters to achieve the best performance for your application. See [Chapter 5, "Optimizing Apache JServ"](#) and [Chapter 4, "Optimizing HTTP Server Performance"](#).
- **Performance monitoring:** Determining what hardware resources are being used by your application and what **response time** your users are experiencing. See [Chapter 2, "Monitoring Your Web Server"](#).
- **Troubleshooting:** Diagnosing why an application is using excessive hardware resources, or why the **response time** exceeds the desired limit.

See Also:

- [Chapter 3, "Sizing and Configuration"](#), for more information on sizing and configuration
- [Chapter 4, "Optimizing HTTP Server Performance"](#), and [Chapter 5, "Optimizing Apache JServ"](#), for more information on parameter tuning
- [Chapter 2, "Monitoring Your Web Server"](#), for more information on performance monitoring

Architecture

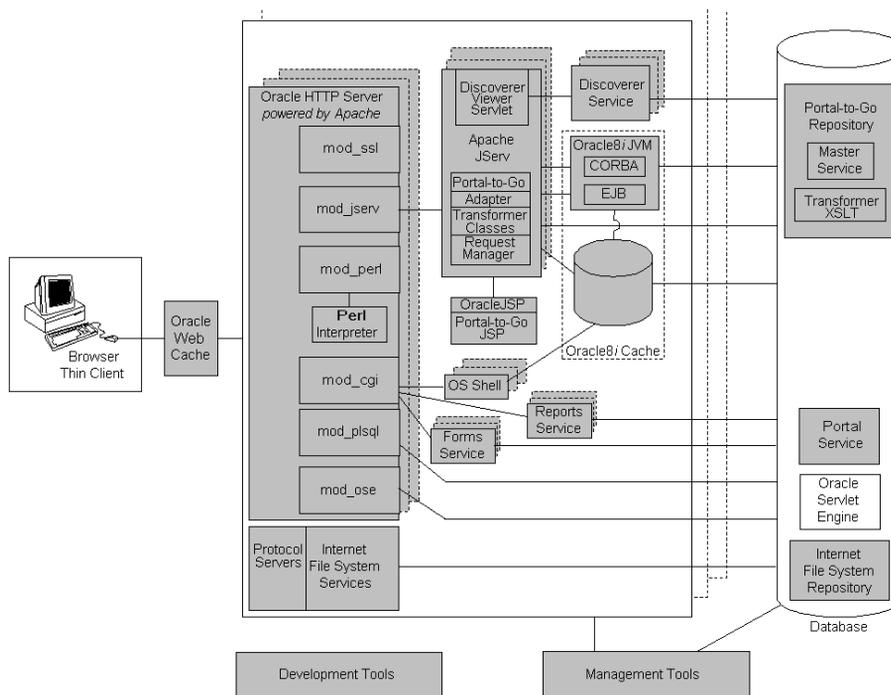
Figure 1-6 shows the architecture of Oracle9i Application Server.

This guide addresses the performance and configuration of these components:

- Oracle HTTP Server *powered by Apache*
- Apache JServ
- OracleJSP

See Also: The Oracle9i Application Server *Overview Guide* for a list of publications that describe other components.

Figure 1-6 Oracle9i Application Server architecture



Monitoring Your Web Server

This chapter describes utilities and processes you can use to gather performance information from your system. This information helps you to determine the best use of your resources.

This chapter contains the following sections:

- [Monitoring Processor Use](#)
- [Monitoring Network Traffic](#)
- [Monitoring the Web Server](#)
- [Monitoring JServ Processes](#)

Monitoring Processor Use

In order to solve performance problems or tune your system, it is useful to know how the CPU is spending its time (for example, scheduling tasks or performing work). You can determine process utilization by gathering CPU statistics, or evaluate system **scalability** by adding users and increasing the system workload, then monitoring paging, swapping, and CPU utilization. Use utilities such as `sar` (System Activity Reporter) and `mpstat` to monitor process use.

Using the sar Utility

You can use `sar` to sample cumulative activity counters in the operating system at specified intervals.

Report CPU Utilization

To determine process use, use the `sar` command. In the example below, the command options `-u`, `3` and `10` tell `sar` to report CPU utilization in three second intervals ten times.

```
prompt>sar -u 3 10

SunOS jpond-sun 5.6 Generic_105181-19 sun4u    10/12/00

12:17:04   %usr   %sys   %wio   %idle
12:17:07     1     0     0     99
12:17:10     0     2     0     97
12:17:13     3     0     0     96
12:17:16    20     6     1     72
12:17:19    15    14     0     71
12:17:22     6     5     0     89
12:17:25     1     0     0     99
12:17:28    16    13     6     65
12:17:31     1     2     1     96
12:17:34     0     1     0     99

Average     6     4     1     88
```

The statistics above show that the CPU was **88%** idle for the given time interval. If your performance criteria specify that CPU usage must be below a certain percentage, then you can use `sar` to sample usage at a chosen interval during peak load times.

The `sar` command (-u option) provides the following statistics:

Table 2-1 CPU statistics, as reported by the `sar` utility

CPU Statistics	Description
%usr	percentage of time in which the processor is running in user mode
%sys	percentage of processes running in system time
%wio	percentage of time the processor spends waiting on I/O requests
%idle	percentage of time that the processor is idle

Using the `mpstat` Utility

The `mpstat` utility is similar to `sar` in that its first argument is the polling interval time in seconds. The second argument to `sarsar` is the number of iterations.

The `mpstat` command:

```
$ mpstat 1 3
```

reports three processor statistics in one second intervals. For example:

```
$ mpstat 1 3
CPU minf mjf xcal intr ithr csw icsw migr smtx srw syscl usr sys wt idl
  0   1   0   0  268  64 148  11   0   0   0   33   3   5   0  92
  0   5   0   0  250  49 157  13   0   1   0  357   2   0   0  98
  0   0   0   0  247  47 134   8   0   0   0  326   0   0   0 100
```

The `mpstat` utility reports the statistics per processor, as shown in [Table 2-2](#).

Table 2-2 CPU statistics, as reported by the `mpstat` utility

Statistic	Description
CPU	processor ID
minf	number of minor faults
mjf	number of major faults
xcal	number of inter-processor cross calls
Intr	number of interrupts

Table 2-2 CPU statistics, as reported by the `mpstat` utility

Statistic	Description
<code>ithr</code>	number of interrupts as threads
<code>csw</code>	number of context switches
<code>icsw</code>	number of involuntary context switches
<code>migr</code>	number of thread migrations to another processor
<code>smtx</code>	number of spins for a mutex (mutual exclusion) lock, which means the lock was not obtained on the first attempt
<code>srw</code>	number of spins on reader-writer lock, which means the lock was not obtained on the first attempt
<code>syscl</code>	number of system calls
<code>usr</code>	percentage of time the processor spent in user mode
<code>sys</code>	percentage of time that the processor spent in system time
<code>wt</code>	percentage of time that the processor spent in wait time (waiting on an event)
<code>idl</code>	percentage of time that the processor spent in idle time

Monitoring Network Traffic

You can use network monitoring tools, such as `snoop` on Solaris or `Network Monitor` on Windows NT, to verify the status of a request as it is being transmitted across the network.

Using the `snoop` Utility

Following are examples of how you can use the `snoop` utility to examine network packets. Using `snoop` in conjunction with `netstat` provides a good picture of network activity.

You can use the `snoop` command with arguments to specify how you want to capture and/or save network packets.

Command	Result
<code>snoop</code>	Captures and displays all packets as they are received.

Command	Result
<code>snoop Athena</code>	Captures and displays all incoming and outgoing packets from host Athena.
<code>snoop -O Gods Athena Zeus</code>	Captures all incoming and outgoing packets between hosts Athena and Zeus, and saves them to a file named <code>Gods</code> .

You can use different command options to view packets captured in a file. For example, the command below displays the contents of the `Gods` file with timestamps relative to the first packet displayed.

```
prompt>snoop -i Gods -t r | more
```

Below is an example of using `snoop` to diagnose a suspected problem related to the `FIN_WAIT_2` state:

```
prompt>snoop -i Gods | grep FIN
```

The first column of the output contains the packet numbers; you can get detailed information about a packet by typing:

```
prompt>snoop -i Gods -v -p<packet number>
```

Monitoring the Web Server

Monitoring activity on the system is essential to performance tuning. The Oracle HTTP Server provides server side status information, including current server statistics, via the `mod_status` module. To obtain these server status reports, you must configure the web server as described in the following sections.

Using the `mod_status` Utility to Monitor the Web Server

To enable monitoring, edit the `httpd.conf` file to replace `your_domain.com` with the hostname of the computer from which you want to monitor.

```
<Location /server-status>
    SetHandler server-status
    Order deny, allow
    Deny from all
    Allow from your_domain.com
</Location>
```

Ensure that the `ExtendedStatus` directive is set to `On`, so that the maximum amount of information is displayed.

When you allow access from all domains, instead of just `your_domain.com`, you can monitor the server from machines outside of your domain, but be aware of the security implications of this: your server status is accessible from any site. It is probably best to specify the domain(s) from which you want to monitor your system.

With monitoring enabled, you can view current statistics from `http://hostname:port/server-status` where `hostname:port` is the hostname and port you want to monitor. These statistics help you to gain insight on how busy your system is.

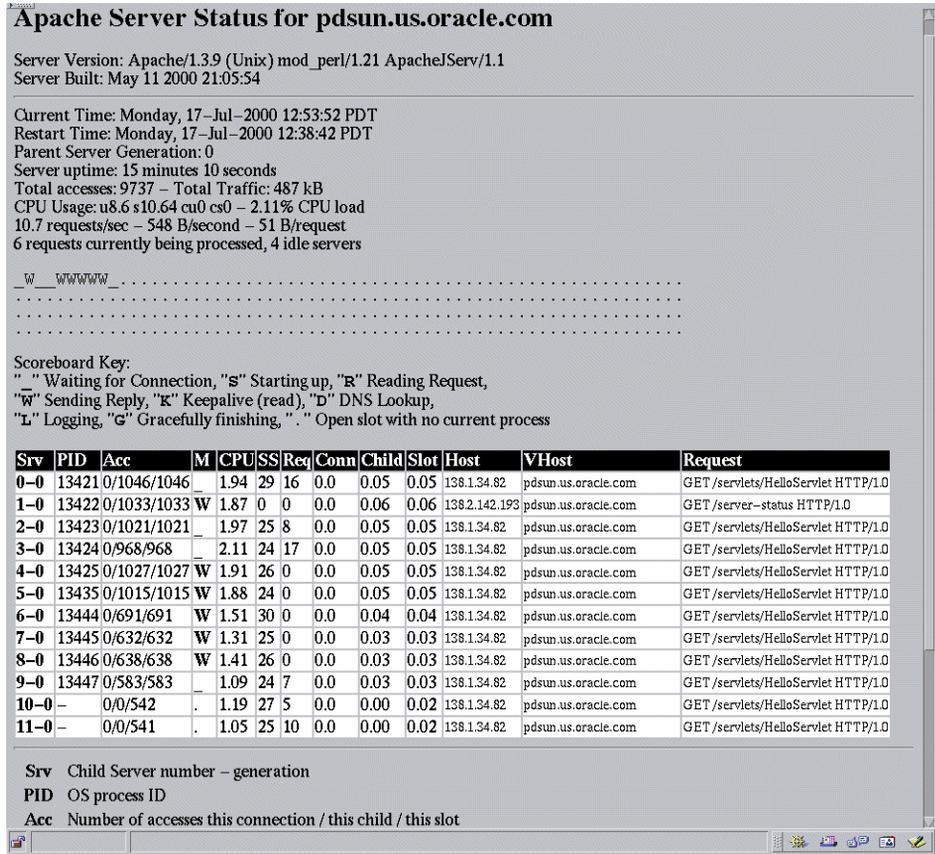
The display includes:

- Hostname for which status is displayed
- Server version
- Date server was built
- Current time, restart time, uptime
- Number of requests currently being processed
- Number of httpd processes serving requests
- Number of idle httpd processes

- Current server state (e.g., waiting for connection, reading request, sending reply, etc).

Figure 2-1 is a screen capture of a server status page with `ExtendedStatus` turned on.

Figure 2-1 Server status page



Interpreting Server Status Information

The display (with `ExtendedStatus` enabled) shows that 6 requests are being processed and four servers are idle. You can determine what stage of processing each server is in from the value in the M (Mode column). In Figure 2-1, 6 servers are sending replies and 4 servers are waiting for connections.

If your system has poor **response times**, or you suspect that httpd processes have stopped responding, look at the Req (request) column. It shows the number of milliseconds required to process the most recent request. Check to see if this number is greater than the time expected to service the request. If, after a request has been completed, there is a W in the M (mode) column for the process, the process is probably not responding.

Another situation that is important to monitor is that of the system being CPU bound, where CPU utilization is around 90%. The server status page displays CPU usage and the number of processes spawned. If the system is approaching the httpd process limit (the `MaxClients` directive's setting in `httpd.conf`), performance is poor, and the processes are all always busy, you may need to change your `MaxClients` setting. See "[Configuring the MaxClients Parameter](#)" on page 4-6.

Customizing the Server Status display

[Figure 2-1](#) is a snapshot of a server for a moment in time. You can get updated server statistics at any interval you choose by including the `refresh` parameter in the server-status URL:

```
http://servername:port/server-status?refresh=x
```

where `servername:port` is the name of the server and port number you are monitoring, and `x` is an integer representing the number of seconds after which the data is refreshed. For example, specify `refresh=3` to update statistics every 3 seconds.

You may also find it useful to have the statistics displayed in a machine-readable format, for processing in a data analysis or spreadsheet program. To do this, add `auto` to the end of the URL, as shown below:

```
http://servername:port/server-status?auto
```

Figure 2-2 Server statistics display

```
Total Accesses: 17503
Total kBytes: 850
CPULoad: 2.6664
Uptime: 1256
ReqPerSec: 13.9355
BytesPerSec: 692.994
BytesPerReq: 49.7286
BusyServers: 1
IdleServers: 9
Scoreboard: _W_.....
```

Logging Server Statistics to a File

The Apache Group provides a Perl script, `log_server_status`, to automate server monitoring. It is included in the `$ORACLE_HOME/Apache/Apache/bin/` directory.

The script is designed to be run by `cron` (or an equivalent daemon that executes commands at intervals). To use the script, you must modify the following configuration variables in the script as shown in [Table 2-3](#).

Table 2-3 *Log status script variables*

Variable	Value
<code>\$wherelog</code>	The pathname of the log file location, for example: <code>/private/admin/logs/</code> The script creates a file name, such as: 20010945.
<code>\$port</code>	Port number of the server to monitor. The default is 80.
<code>\$server</code>	The server host name. The default is <code>localhost</code> .
<code>\$request</code>	The server status request with the <code>auto</code> parameter as entered in the browser, for example: <code>http://servername:port/server-status?auto</code>

Enabling server status is very useful if an `httpd` process is not responding, and you need to identify that process. Operating system utilities such as `ps`, `top`, or `pmap` do not identify which process is not responding.

Monitoring JServ Processes

After you start the Oracle9i Application Server, you can check to ensure that all JServ processes have started normally. If performance is degraded during operation, you can quickly determine if this is because JServ processes have terminated by looking at the Status column (each configured process has a status of Up or Down).

1. Remove the comments in the JServ status handler section of the `jserv.conf` file to enable monitoring and specify the host(s) that can access JServ status (the default is `localhost`). Be aware of security implications when selecting the hosts that will be allowed to access status information on your system.

```
<Location /jserv/>
SetHandler jserv-status
    order deny, allow
    deny from all
    allow from hostname_1.com
    allow from hostname_2.com
</Location>
```

2. Type the following into your browser:

```
http://hostname:port/jserv/
```

The port must be the port on which the web server listens (found in the `httpd.conf` file). Always include the trailing slash (/) in this URL. A “not found” error occurs if you omit the trailing slash.

A Configured Hosts column displays links to hosts.

3. Click the host to monitor.

The JServ status information for the host displays as shown in [Figure 2-3](#).

Figure 2-3 JServ status display

Note: The JServ status monitor shows all of the JServ processes that are configured in the `jserv.conf` file, but not all of these may have been started, or any of them could be terminated. For example, [Figure 2-3](#) shows four processes, but only two have a Status of Up (indicating that the process is able to service requests).

Parameter	Value
Server Name	pdsun-perf.us.oracle.com
ApJServManual	TRUE (STANDALONE OPERATION)
ApJServProperties	/private/oracle/app/oracle/ias/Apache/Jserv/etc/jserv.properties (IGNORED)
ApJServDefaultProtocol	ajpv 12 (PORT 8007)
ApJServDefaultHost	localhost (ADDR 127.0.0.1)
ApJServDefaultPort	8007
ApJServLogFile	/private/oracle/app/oracle/ias/Apache/Jserv/logs/jserv.log (DESCRIPTOR 7)
ApJServMountCopy	TRUE
ApJServShm File	/private/oracle/app/oracle/ias/Apache/Apache/logs/jserv_shm

MountPoint	Server	Protocol	Host	Port	Zone	Status
/servlets/	pdsun.us.oracle.com	balance	set1		root	
	JServ1 weight=1	ajpv 12	127.0.0.1 (ADDR 127.0.0.1)	8007	"	JS1's current shm state: Up (+) <input type="button" value="test"/> change to: <input type="button" value="choose"/> <input type="button" value="apply"/>
	JServ2 weight=1	ajpv 12	127.0.0.1 (ADDR 127.0.0.1)	8008	"	JS2's current shm state: Up (+) <input type="button" value="test"/> change to: <input type="button" value="choose"/> <input type="button" value="apply"/>
	JServ3 weight=1	ajpv 12	127.0.0.1 (ADDR 127.0.0.1)	8009	"	JS3's current shm state: Down (-) <input type="button" value="test"/> change to: <input type="button" value="choose"/> <input type="button" value="apply"/>
	JServ4 weight=1	ajpv 12	127.0.0.1 (ADDR 127.0.0.1)	8010	"	JS4's current shm state: Down (-) <input type="button" value="test"/> change to: <input type="button" value="choose"/> <input type="button" value="apply"/>

The Status column shows the current shared memory (shm) state of each process.

Note: The Status column is populated only for processes that are started in manual mode. It is not populated for a single process started in automatic mode.

The symbols that appear in parentheses after the word Up or Down have the following meanings:

Symbol	Meaning
+	The process is running.
-	The process is stopped.
x	The process was terminated in a harsh shutdown. (existing requests were not handled before the process was terminated).
/	The process was terminated in a graceful shutdown (existing requests were handled before the process was terminated).

Sizing and Configuration

This chapter provides guidelines for sizing and configuration which can help you meet performance goals. It also discusses performance factors such as CPU and memory consumption.

This chapter contains the following sections:

- [Sizing your Hardware and Resources](#)
- [Understanding Concurrent Users and User Population](#)
- [Determining CPU Requirements](#)
- [Determining Memory Requirements](#)

Sizing your Hardware and Resources

In addition to the minimum installation recommendations, your hardware resources need to be adequate for the requirements of your specific applications. To avoid hardware-related performance bottlenecks, each hardware component should operate at no more than 80% of capacity. See ["Using the sar Utility"](#) on page 2-2 for information on measuring CPU utilization.

Processor and memory resources in particular should be allocated generously, for the maximum user load expected.

See Also: ["Using the sar Utility"](#) on page 2-2 for information on measuring CPU utilization.

Understanding Concurrent Users and User Population

The amount of hardware resources required varies based on the application. A common mistake is to use resource estimates that do not incorporate user **think time** and network latencies. In sizing applications, you must have some idea of the relationship between the number of potential users and the number of concurrent users. This is determined by the **think time** and the average **response time** for your application.

To determine memory requirements, you also need to consider the number of concurrent executing users (not the total user population) times the cost per user.

Note: `TMaxClients` setting in your `httpd.conf` file limits the number of concurrently executing users. See ["Configuring the MaxClients Parameter"](#) on page 4-6 for more information.

[Table 3-1](#) provides an example of the impact of **think time** and **service time** on the **concurrency** and resulting performance of a system.

Table 3-1 *Concurrent executing users*

User population ¹	Think time (sec) ²	Service time (sec) ³	Range of concurrent users ⁴	Average response Time (sec) ⁵	Requests per second (throughput) ⁶	CPU utilization (%) ⁷
100	0	0.3	100	5.2	19	99
100	1	0.3	65-100	4.2	19	99
100	10	0.3	0-32	0.9	9	48
100	10	0.6	0-53	2.9	8	80

¹ User population - total users.

² **Think time** - the time the user is not engaged in actual use of the processor (the time between requests).

³ **Service time** (seconds) - elapsed time to complete the operation measured for a single user.

⁴ Range of concurrent users - the number of users measured on the server, taken in snapshots from the `server-status` display (requests currently being processed). See ["Using the mod_status Utility to Monitor the Web Server"](#) on page 2-6 for information on `server-status`.

⁵ Average **response time** - **response time** measured at the client under load.

⁶ Requests per second (**throughput**) - number of requests processed.

⁷ CPU utilization - average total CPU utilization as a percentage.

Determining CPU Requirements

For most applications, the majority of the CPU utilization is spent in processing the application's code. The CPU requirement of any application depends on its complexity and workload, as shown in [Table 3-2](#).

You will need to monitor the CPU requirements of applications throughout the development cycle. See [Chapter 2, "Monitoring Your Web Server"](#) for information on how to do this.

Table 3-2 *Application CPU requirements on a336 MHz Sun SPARC processor*

Application	CPU requirement (per request)
Static page, 20KB	5 msec
Simple servlet, JDK (Java Developer's Kit) release 1.2	20 msec
Simple servlet, JDK release 1.1.8	40 msec

Table 3–2 Application CPU requirements on a336 MHz Sun SPARC processor

Application	CPU requirement (per request)
Medium application	100-200 msec
Complex application	400-600 msec

Secure Sockets Layer Impact on CPU Requirements

Secure Sockets Layer (SSL) is a protocol used for transmitting documents securely over the Internet. URLs for Web pages that require an SSL connection begin with `https` instead of `http`.

Establishing an SSL connection is costly in terms of **response time** and CPU utilization. For example, a request with a **response time** of 0.5 seconds without SSL generated a **response time** of 1.7 seconds with SSL (measured on an internal 100 Mbps network). Most of the performance cost in using SSL is in establishing the connection (approximately 125 ms of CPU time per connection on a 336 Mhz processor).

The high connection cost is incurred for the first connection in a client's SSL session, because the HTTP Server can cache the SSL session information, reducing the overhead for subsequent connections. For more information, see "[Enabling SSL Session Caching](#)" on page 4-7.

Determining Memory Requirements

This section discusses the following memory requirements:

- [Determining Memory Requirements for Operating System Software](#)
- [Determining Memory Requirements for the Oracle HTTP Server](#)
- [JServ Memory Requirements](#)
- [Determining Java Heap Size](#)
- [Determining Memory Requirements for Servlets and OracleJSP pages](#)
- [Determining the Number of JServ Processes per CPU](#)

Determining Memory Requirements for Operating System Software

In an idle system with memory resources freely available, your operating system statistics may indicate that the resident memory usage is close to the virtual size. As users place more load on the system, the operating system reclaims unneeded memory from these processes, and the amount of resident memory they consume decreases. If you are monitoring your own system, take snapshots of processes at varying usage levels.

Refer to your operating system hardware and software documentation for more information on measuring and tuning operating system memory usage. You can monitor memory usage and processor statistics with standard operating system tools. See [Chapter 2, "Monitoring Your Web Server"](#) for more information.

Sun recommends reserving 15% of the overall real memory on the system for the kernel and other system overhead.

Determining Memory Requirements for the Oracle HTTP Server

In a series of tests of listener memory usage, each HTTP listener used (at startup) approximately 400K of resident memory. This size increased by 500-600K per process when the listener was active. When it was dormant, the operating system reduced the listener's memory usage back to the startup size.

Using standard operating system tools, you can examine resident memory sizes. If you look at a listener process, you will see that it is larger than the figure above because the displayed size includes shared memory.

JServ Memory Requirements

A JServ process using JDK release 1.2 requires 12-15 MB at startup. Using JDK release 1.1.8, it requires 10 MB.

Determining Java Heap Size

For JDK release 1.1.8, the default maximum heap size is 16MB. For JDK release 1.2, it is 24MB.

To maximize performance, set the maximum heap size to accommodate application requirements. To determine how much Java heap you need, include calls in your program to the `Runtime.getRuntime().totalMemory()` and `Runtime.getRuntime().freeMemory()` methods in the `java.lang` package. Subtract free memory from total memory; the difference is the amount of heap that the application consumed.

Suppose you determine that you need 128MB of heap. To change the heap size, you would set the maximum Java heap size in the `jserv.properties` file for automatic mode:

```
wrapper.bin.parameters=-mx128m
```

In manual mode, if more than one JServ process is running, the heap size must be set on the command line or in the startup script for each JServ process.

When a JServ process exceeds its maximum heap size, the process terminates. In automatic mode, a new process is started, but performance is degraded significantly. In manual mode, a terminated process will not be restarted, so ensure that the heap size is sufficient.

Note: The process size reported by utilities such as `top` or `ps` will be larger than the maximum heap size, because private memory is added to the maximum heap size.

Determining Memory Requirements for Servlets and OracleJSP pages

OracleJSP pages (Oracle's implementation of Sun's JavaServer Pages) and servlets require different amounts of memory, depending on the release of the JDK used. [Table 3-3](#) compares memory requirements for a simple servlet and an Oracle JSP page under load with 10-30 active threads. The servlet did not use sessions. The OracleJSP page had sessions on (the default).

Table 3-3 *Servlet and OracleJSP pages memory*

Component	JDK release 1.1.8	JDK release 1.2
Servlet	10MB	24MB
OracleJSP page	10MB	32MB

The amount of memory needed depends on whether sessions are used; a session consumes about 0.5KB. For maximum performance, if sessions are not being used, turn them off in the OracleJSP application as follows:

```
<%@ page session="false" %>
<html><body>
HelloWorld
</body></html>
```

As a starting point, figure that each active user consumes at least 150 KB to 200 KB for Java applications, plus the size of the server processes. For Java applications, the base process size is approximately 12-15 MB.

An application's memory needs also depend on its size, the amount of data cached, and other factors.

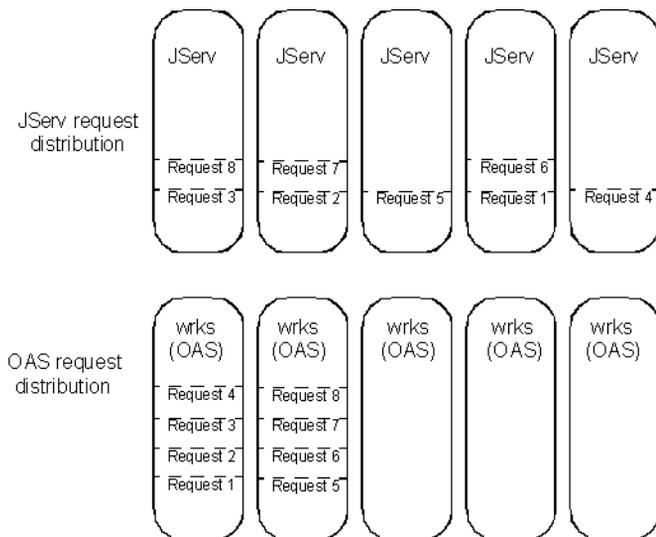
See Also: The *OracleJSP Developer's Guide and Reference* in the Oracle9i Application Server documentation library for more information on OracleJSP pages.

Determining the Number of JServ Processes per CPU

Oracle recommends two JServ processes per CPU as a starting point. The default thread setting (`security.maxConnections=50`) in the JServ configuration file (`jserv.conf`) is also a good starting point.

If your application code performs a lot of synchronization, or creates many new Java objects, then you should consider increasing the number of JServ processes, while limiting the number of threads per process to between 10 and 20. In this way you avoid increased queuing and processing required for object synchronization in the JVM (Java virtual machine). This is because the `httpd (mod_jserv)` process sends incoming requests to the JServ processes in a distributed fashion. (Readers familiar with the Oracle Application Server will recall that requests are sent to a servlet engine until its thread limit is reached, and subsequent requests are sent to the next servlet engine.) This difference is illustrated in [Figure 3-1](#).

Figure 3–1 Request distribution



See Also: ["How to Perform Load Balancing"](#) on page 5-4 for instructions on changing parameters in the configuration files, and details on how requests are distributed among the available JServ engines.

Optimizing HTTP Server Performance

This chapter provides information on improving the Oracle HTTP Server's performance, including tuning TCP parameters, the effects of changing the `MaxClients` parameter, SSL caching, and the performance impacts of logging.

This chapter contains the following sections:

- [Tuning TCP Parameters](#)
- [Configuring the MaxClients Parameter](#)
- [Enabling SSL Session Caching](#)
- [Understanding Performance Implications of Logging](#)
- [Benefits of the HTTP/1.1 Protocol](#)
- [Differences between Apache Releases with Respect to Performance](#)

Tuning TCP Parameters

Correctly tuned TCP parameters can improve performance dramatically. This section contains recommendations for TCP tuning and a brief explanation of each parameter.

[Table 4-1](#) contains recommended TCP parameter settings and references to discussions of each parameter in this section.

Table 4-1 Recommended TCP parameter settings

Parameter	Setting	Comments
<code>tcp_conn_hash_size</code>	32768	See "Increasing TCP Connection Table Access Speed" on page 4-3.
<code>tcp_close_wait_interval</code>	60000	Parameter name in Solaris release 2.6. See "Specifying Retention time for Connection Table entries" on page 4-3.
<code>tcp_time_wait_interval</code>	60000	Parameter name in Solaris release 2.7. See "Specifying Retention time for Connection Table entries" on page 4-3.
<code>tcp_conn_req_max_q</code>	1024	See "Increasing the Handshake Queue Length" on page 4-4.
<code>tcp_conn_req_max_q0</code>	1024	See "Increasing the Handshake Queue Length" on page 4-4.
<code>tcp_slow_start_initial</code>	2	See "Changing the Data Transmission Rate" on page 4-4.
<code>tcp_xmit_hiwat</code>	32768	See "Changing the Data Transfer Window Size" on page 4-5.
<code>tcp_rcv_hiwat</code>	32768	See "Changing the Data Transfer Window Size" on page 4-5.

Setting TCP parameters

To set the connection table **hash** parameter, you must add the following line to your `/etc/system` file, and then restart the system:

```
prompt>set tcp:tcp_conn_hash_size=32768
```

A sample script, `tcpset.sh`, that changes TCP parameters to the settings recommended here, is included in the `$ORACLE_HOME/Apache/Apache/bin/` directory.

If your system is restarted after you run the script, the default settings will be restored and you will have to run the script again. To make the settings permanent, enter them in your system startup file.

Increasing TCP Connection Table Access Speed

If you have a large user population, you should increase the **hash** size for the TCP connection table. The **hash** size is the number of **hash** buckets used to store the connection data. If the buckets are very full, it takes more time to find a connection. Increasing the **hash** size will reduce the connection lookup time, but increases memory consumption.

Suppose your system performs 100 connections per second. If you set `tcp_close_wait_interval` to 60000, then there will be about 6000 entries in your TCP connection table at any time. Increasing your **hash** size to 2048 or 4096 will improve performance significantly.

On a system servicing 300 connections per second, changing the **hash** size from the default of 256 to a number close to the number of connection table entries decreases the average round trip time by three to four seconds. The maximum **hash** size is 262144. Ensure that you increase memory as needed.

To set the `tcp_conn_hash_size`, add the line shown below to your `/etc/system` file. The parameter will take effect when the system is restarted.

```
set tcp:tcp_conn_hash_size=32768
```

Specifying Retention time for Connection Table entries

The TCP connection table maintains data associated with connections. The server maintains an entry in this table for some time after a connection is closed, so that it can identify and properly dispose of any leftover incoming packets from the client.

Access speed to this table impacts performance; the access speed depends on the number of entries in the table, and on its hash size. The number of entries in the table depends on the rate of incoming requests, and the lifetime of each connection.

You can control the length of time that TCP connection table entries are maintained with the `tcp_close_wait_interval` parameter (renamed `tcp_time_wait_interval` on Solaris release 2.7). This parameter is commonly set to 60,000 msec. Use the following command to set it (note the difference in parameter name for Solaris release 2.6 and 2.7).

In Solaris release 2.6:

```
prompt>/usr/sbin/ndd -set /dev/tcp tcp_close_wait_interval 60000
```

In Solaris release 2.7:

```
prompt>/usr/sbin/ndd -set /dev/tcp tcp_time_wait_interval 60000
```

Note: If your user population is widely dispersed (with respect to Internet topology), you may want to set this parameter to a higher value. You can improve access time to the TCP connection table with the `tcp_conn_hash_size` parameter.

Increasing the Handshake Queue Length

During the TCP connection handshake, the server, after receiving a request from a client, sends a reply, and waits to hear back from the client. The client responds to the server's message and the handshake is complete. Upon receiving the first request from the client, the server makes an entry in the listen queue. After the client responds to the server's message, it is moved to the queue for messages with completed handshakes. The second queue makes it possible for the server to continue servicing requests for which the handshake has been completed.

The maximum length of the queue for incomplete handshakes is governed by `tcp_conn_req_max_q0`, which by default is 1024. The maximum length of the queue for requests with completed handshakes is defined by `tcp_conn_req_max_q` (default is 128).

On most web servers, the defaults will be sufficient, but if you have more than 1024 concurrent users, these settings may be too low. In that case, connections will be dropped in the handshake state because the queues are full. You can determine whether this is a problem on your system by inspecting the values for `tcpListenDrop`, `tcpListenDropQ0`, and `tcpHalfOpenDrop` with `netstat -s`. If either of the first two values are nonzero, you should increase the maximums.

The defaults are probably sufficient, but Oracle recommends that you increase the value of `tcp_conn_req_max_q` to 1024. You can set these parameters with:

```
prompt>/usr/sbin/ndd -set /dev/tcp tcp_conn_req_max_q 1024
prompt>/usr/sbin/ndd -set /dev/tcp tcp_conn_req_max_q0 1024
```

Changing the Data Transmission Rate

Typically, all packets in a data transfer are sent at once. TCP implements a slow starting data transfer to prevent overloading a busy segment of the Internet. With

slow start, one packet is sent, an acknowledgment is received, then two packets are sent. The number sent to the server continues to be doubled after each acknowledgment, until the TCP transfer window limits are reached.

Some versions of Microsoft Windows (including NT 4.0 and 95) do not acknowledge receipt of a single packet when a connection is initiated, but if two packets are received, an acknowledgment is sent immediately. Because Solaris sends only one packet when initiating a connection (per the TCP standard), this can increase the connection startup time. This is especially apparent on fast local networks, where the **latency** is expected to be low.

To configure Solaris to start with two packets when initiating a data transfer, use the following command:

```
prompt>/usr/sbin/ndd -set /dev/tcp tcp_slow_start_initial 2
```

Changing the Data Transfer Window Size

The size of the TCP transfer windows for sending and receiving data determine how much data can be sent without waiting for an acknowledgment. The default window size is 8192 bytes. Unless your system is memory constrained, these windows should be increased to the maximum size of 32768. This can speed up large data transfers significantly. Use these commands to enlarge the window:

```
prompt>/usr/sbin/ndd -set /dev/tcp tcp_xmit_hiwat 32768
```

```
prompt>/usr/sbin/ndd -set /dev/tcp tcp_rcv_hiwat 32768
```

Because the client typically receives the bulk of the data, it would help to enlarge the TCP receive windows on end users' systems.

Configuring the MaxClients Parameter

The `MaxClients` directive limits the number of clients that can simultaneously connect to your web server, and thus the number of `httpd` processes. You can configure this parameter in the `httpd.conf` file up to a maximum of 1024 in Oracle9i Application Server release 1.0.2 (in the previous release, the maximum was 256). The default is 150, which should be adequate for most uses. If the `MaxClients` setting is too low, and the limit is reached, clients will be unable to connect.

Our tests of static page requests (average size 20K) on a 2 processor, 168 MHz Sun UltraSPARC on a 100 Mbps network showed that:

- The default `MaxClients` setting of 150 was sufficient to saturate the network.
- Approximately 60 `httpd` processes were required to support 300 users (no **think time**).

On the system described above, and on 4 and 6-processor, 336 MHz systems, there was no significant performance improvement in increasing the `MaxClients` setting from 150 to 256, based on static page and servlet tests with up to 1000 users.

Increasing `MaxClients` when system resources are saturated does not improve performance. When there are no `httpd` processes available, connection requests are queued in the TCP/IP system until a process becomes available, and eventually clients terminate connections.

Note: If you are using persistent connections, you may require more concurrent `httpd` server processes. See "[How Persistent Connections Can Reduce httpd Process Availability](#)" on page 4-10 for a discussion of the relationship between persistent connections and the number of server processes.

For dynamic requests, if the system is heavily loaded, it might be better to allow the requests to queue in the network (thereby keeping the load on the system manageable). The question for the system administrator is whether a timeout error and retry is better than a long **response time**. In this case, the `MaxClients` setting could be reduced, to act as a throttle on the number of concurrent requests on the server.

Enabling SSL Session Caching

The Oracle HTTP server caches a client's SSL session information by default. With session caching, only the first connection to the server incurs high **latency**. For example, in a simple test to connect and disconnect to an SSL-enabled server, the elapsed time for 5 connections was 11.4 seconds without SSL session caching. With SSL session caching enabled, the elapsed time for 5 round trips was 1.9 seconds.

The `SSLSessionCacheTimeout` directive in `httpd.conf` determines how long the server keeps a session alive (the default is 300 seconds). The session information is kept in a file. You can specify where to keep the session information using the `SSLSessionCache` directive; the default location is the `$ORACLE_HOME/Apache/Apache/logs/` directory. The file can be used by multiple Oracle HTTP Server processes.

The duration of an SSL session is unrelated to the use of HTTP persistent connections.

Understanding Performance Implications of Logging

This section discusses types of logging, log levels, and the performance implications of using them.

Access Logging

For static page requests, access logging of the default fields results in a 2-3% performance cost.

HostNameLookups

By default, the `HostNameLookups` directive is set to `Off`. The server writes the IP addresses of incoming requests to the log files. When `HostNameLookups` is set to `on`, the server queries the DNS system on the Internet to find the host name associated with the IP address of each request, then writes the host names to the log.

Performance degraded by about 3% (best case) in Oracle in-house tests with `HostNameLookups` set to `on`. Depending on the server load and the network connectivity to your DNS server, the performance cost of the DNS lookup could be high. Unless you really need to have host names in your logs in real time, it is best to log IP addresses. You can resolve IP addresses to host names off-line, with the `logresolve` utility (found in the `$ORACLE_HOME/Apache/Apache/bin/` directory).

Error logging

The server notes unusual activity in an error log. The `ErrorLog` and `LogLevel` directives identify the log file and the level of detail of the messages recorded. The default level is `warn`. There was no difference in static page performance on a loaded system between the `warn`, `info`, and `debug` levels.

Benefits of the HTTP/1.1 Protocol

The Oracle HTTP server can use HTTP/1.1. Netscape Navigator release 4.0 still uses HTTP/1.0, with some 1.1 features, such as persistent connections. Internet Explorer uses HTTP/1.1. The performance benefit of persistent connections comes from reducing the overhead of establishing and tearing down a connection for each request. A persistent connection accepts multiple requests from a user.

For a small static page request, the connection **latency** can equal or exceed the response **latency** (the time to fulfill the request after the connection is established), so using persistent connections can result in major performance gains.

Supporting Persistent Connections

If your users' browsers support persistent connections (the default behavior of HTTP/1.1), you can support them on the server using the `KeepAlive` directives in the Oracle HTTP Server. (Some browsers that do not support all HTTP/1.1 features do support persistent connections; for example, recent versions of Netscape.)

How Persistent Connections Improve Response Times

Persistent connections can improve total **response time** for a web interaction that involves multiple HTTP requests, because the delay of setting up a connection only happens once.

Consider the total time required, without persistent connections, for a client to retrieve a web page with three images from the server.

Activity	Seconds
Establish connection	1
Produce and send the text portion of the page	5
Establish connection	1
Transfer first image file	2

Establish connection	1
Transfer second image file	2
Establish connection	1
Transfer third image file	2
Total	15

With persistent connections, the **response time** for the same request is reduced:

Activity	Seconds
Establish connection	1
Produce and send the text portion of the page	5
Transfer first image file	2
Transfer second image file	2
Transfer third image file	2
Total	12

This is a 20% reduction in **service time**.

How Persistent Connections Reduce Server Workload

Another benefit of persistent connections is reduction of the work load on the server. Because the server need not repeat the work to set up the connection with a client, it is free to perform other work.

For a very inexpensive servlet (Hello World), the CPU ms per request was reduced by approximately 10% when the same client made 4 requests per connection. (The impact would be far less significant for a realistic servlet application that does more work.) In-house tests using an OracleJSP application (`lotto.jsp`, one of the samples that ships with Oracle9i Application Server) and persistent connections showed an improvement of about 20%, with a single user making 5 requests per connection. With an increased number of users (10-100), the performance improvement was less dramatic, but still significant (6% or better).

How Persistent Connections Can Reduce httpd Process Availability

There are some serious drawbacks to using persistent connections with Apache. In particular, because httpd processes are single threaded, one client can keep a process tied up for a significant period of time (the amount of time depends on your `KeepAlive` settings). If you have a large user population, and you set your `KeepAlive` limits too high, clients could be turned away because of insufficient httpd daemons.

The default settings for the `KeepAlive` directives are:

```
KeepAlive on
MaxKeepAliveRequests 100
KeepAliveTimeout 15
```

These settings allow enough requests per connection and time between requests to reap the benefits of the persistent connections, while minimizing the drawbacks. You should consider the size and behavior of your own user population in setting these values on your system. For example, if you have a large user population and the users make small infrequent requests, you may want to reduce the above settings, or even set `KeepAlive` to `off`. If you have a small population of users that return to your site frequently, you may want to increase the settings.

Understanding FIN_WAIT_2 Connection Problems

There is a known problem with some browsers which will leave the server with a TCP connection in the `FIN_WAIT_2` state. If too many connections are left in this state, the system will run out of the memory allocated for storing TCP connections, and stop.

The problem is that when a connection becomes idle, and the server closes upon expiration of the keep alive time limit, the client host may not perform the TCP protocol steps required to complete the closure of the connection. The host, having sent the close request, is left with the connection in the `FIN_WAIT_2` state consuming memory until it gets the appropriate packets back from the client, or until an internal flush occurs. If a connection is left in the `FIN_WAIT_2` state, the httpd process with which the connection is associated is freed to service other requests as indicated, so this problem won't tie up web server processes.

On Solaris, the parameter `tcp_fin_wait_2_flush_interval` dictates the frequency with which these connections will be cleaned up. In general, the default setting is sufficient, and should not be modified unless the system is failing.

Note: The `FIN_WAIT_2` state can also occur due to a system bug unrelated to use of `KeepAlive`. The bug is fixed by the Solaris cluster patch 105181-20.

Differences between Apache Releases with Respect to Performance

The difference between Apache releases 1.3.9 and 1.3.12 was primarily corrected bugs. With static page and servlet performance measurements, there was no performance difference measured between the versions.

Optimizing Apache JServ

This chapter describes the JServ architecture, and discusses ways you can improve JServ performance. It also includes performance information on OracleJSP pages (the Oracle implementation of Sun Microsystems' JavaServer Pages 1.1.)

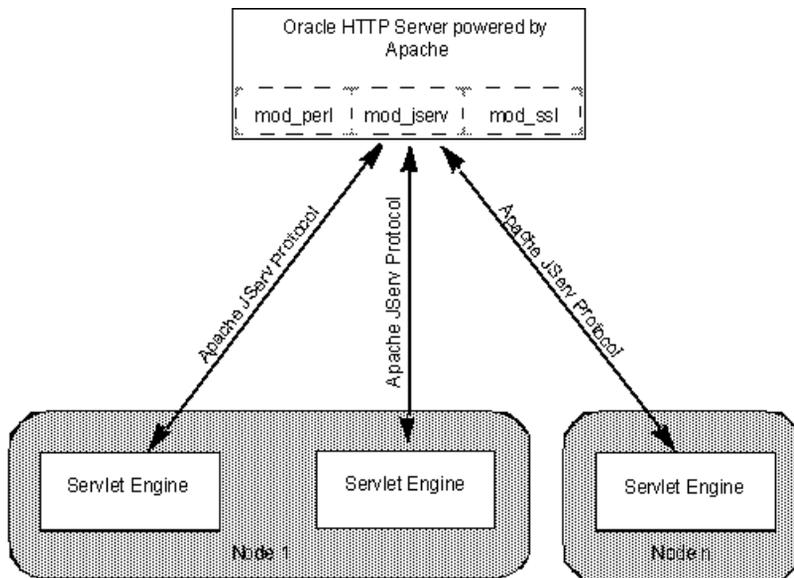
- [Overview of JServ](#)
- [Optimizing Servlet Performance](#)
- [What is OracleJSP?](#)
- [Tuning OracleJSP Pages for Performance](#)

Overview of JServ

Apache JServ is made up of an Apache module called `mod_jserv`, which runs in the `httpd` process, and a servlet engine, which runs in a Java process. `mod_jserv`, which is implemented in C, functions as a dispatcher, routing each servlet request to a JServ process for execution.

The servlet engine runs in its own JVM (Java Virtual Machine) and is solely responsible for parsing the request and generating a response. As [Figure 5-1](#) shows, multiple JServs can service requests. The HTTP server process and the JServ process communicate using the Apache JServ Protocol 1.2.

Figure 5-1 Apache JServ components



Optimizing Servlet Performance

This section discusses strategies for optimizing JServ performance: loading servlets when starting the JVM, and load balancing.

The terms “repository” and “zone” are used in this discussion. Servlets, repositories, and zones are analogous to files, directories and virtual hosts. A servlet is a single unit, a **repository** is a collection of servlets, and a **zone** is a collection of repositories.

Loading Servlet Classes

Apache JServ allows you to load servlet classes when the JVM is started. To do this, put the servlets to load in the `servlets.startup` directive in the servlet zone properties file. When the servlet is loaded, its `init()` method is called. All other servlets (those not listed in `servlets.startup`) are loaded and initialized on first request.

Using this facility increases the start-up time for your JServ process, but improves first-request **latency** for servlets.

Pre-Loading with JSPs

If you are using a JSP as the servlet (your code does not extend `HttpServlet`), you will be unable to use this pre-load option, but you could pre-load the JSP runner by including the `oracle.jsp.jspServlet` in `servlets.startup`.

If the first-request **latency** for your initialization routines is really a performance issue, you can achieve some of the results described above by creating a dummy servlet to call your one-time initialization routines in its `init()` method. You must add the name of the dummy servlet to `servlets.startup`.

Reloading Servlet Classes Automatically

If `autoreload.classes` is set to `true` for a zone (the default), then each time one of that zone’s servlets is requested, every class that has been loaded from a repository in that zone is checked to see if it has been modified. If one of the classes has changed, then all previously loaded classes from the zone’s repositories are unloaded, which means that as the classes are needed, they will be loaded from their class files again.

This is a useful development feature, because you can install new versions or drop in new class files without restarting the server. For optimal performance in production environments, however, you should set both automatic class reloading

parameters to false, since there is a performance cost in checking the repositories on every execution of a servlet. Change these parameters in the zone properties file:

```
autoreload.classes=false
autoreload.file=false
```

How to Perform Load Balancing

It is often beneficial to spread the servlet application load among multiple JServ processes, especially when the application is run on a multiprocessor system or if the servlets and HTTP server are run on separate nodes. Running multiple Apache JServ processes generally results in higher **throughput** and shorter **response time**, even on a single-processor host. (See [Chapter 3, "Sizing and Configuration"](#) for specific recommendations.)

Using Automatic Load Balancing

The `mod_oprocmgr` module shipped with Release 1.0.2.2 enables you to specify a 'auto' setting in `jserv.conf` for the `ApJServManual` directive. `mod_oprocmgr` provides automatic load balancing and death detection for JServ processes. See *Using mod_oprocmgr with mod_jserv* in the Oracle9i Application Server documentation library for instructions on using this module.

Using Manual Load balancing

This section explains how to balance incoming requests between two JServ processes running on the same host as the HTTP server. Examples from the `jserv.properties` files are included with the procedures; substitute your own port numbers and directory locations where needed.

With this method of load balancing, you must start and stop processes manually, because JServ cannot automatically start and stop more than one JServ process. (Sample scripts for starting and stopping the JServ processes and the Oracle HTTP Server are included in the `$ORACLE_HOME/Apache/Apache/bin/` directory.) This means that if a process terminates for any reason, JServ will not restart it. To prevent processes from terminating due to memory shortage, ensure that you have a sufficient maximum heap size set for your JServ processes. See "[Determining Java Heap Size](#)" on page 3-5.

Configuring the JServ processes

Each JServ process in your load balancing scheme must be configured to listen on its own port and to log to its own file. If you have a `jserv.properties` file

containing the parameters needed to run your application, you can duplicate it to create a properties file for each JServ process.

1. Create a properties file for each JServ process.

```
prompt>cp jserv.properties jserv1.properties
prompt>cp jserv.properties jserv2.properties
```

2. Edit *jserv1.properties* as follows:

```
port=8001
log.file=/usr/local/jserv/logs/jserv1.log
```

3. Edit *jserv2.properties* as follows:

```
port=8002
log.file=/usr/local/jserv/logs/jserv2.log
```

If JServ is included in your CLASSPATH, you can start the JServ processes with

Note: If your HTTP server will be running on a different host than the JServ processes, you must also add the IP address of the host running the HTTP server to the *security.allowedAddresses* parameter in each *jserv.properties* file.

these commands:

```
java JServ jserv1.properties
java JServ jserv2.properties
```

To start and stop the processes and the web server, it is convenient to use scripts. Samples are included in the `$ORACLE_HOME/Apache/Apache/bin/` directory (*startJServ.sh* and *stopJServ.sh*).

Modifying *jserv.conf* to distribute the load

1. Set the flag to start processes manually.

```
ApJServManual on
```

2. Indicate where the servlet request is to be sent.

- a. Locate the *ApJServMount* directive.

```
ApJServMount /servlets /root
```

If the user requests

`http://your.server.com/servlets/testServlet`, the `ApJServMount` directive above will execute `testServlet` in the **zone** called `/root`.

- b. Change the **zone** identifier from `/root` to `balance://Jserv_set/root` and then add the `ApJServBalance`, `ApJServHost`, `ApJServRoute` directives for each process sharing the load, as shown below:

```
ApJServMount /servlets balance://Jserv_set/root
ApJServBalance JServ_set JServ1
ApJServBalance JServ_set JServ2 2
ApJServHost JServ1 ajpv12://127.0.0.1:8001
ApJServHost JServ2 ajpv12://127.0.0.1:8002
ApJServRoute JS1 JServ1
ApJServRoute JS2 JServ2
ApJServShmFile /usr/local/apache/logs/jserv_shm
```

- * The `ApJServMount` directive, with `/servlets balance://Jserv_set/root`, now balances requests for servlets in `/servlets` between `JServ1` and `JServ2`.
- * The `ApJServBalance` directive identifies `JServ1` and `JServ2` as the processes that share the load. The '2' following `JServ2` is a weight value. It specifies that twice as many requests will be sent to `JServ2` as would be otherwise, i.e., that `JServ2` will get about 2/3 of all incoming requests. See "[Distribution of JServ Requests](#)" below for details.
- * The `ApJServHost` directive identifies the host and port on which the processes are listening.
- * The `ApJServRoute` directive associates `JServ` processes with sessions. `JServ` uses this information to keep all of a session's requests together in one process. The `JServ` session mechanism sends the process route information back to the user (generally in a cookie). You need only modify it if your application uses sessions.
- * The `ApJServShmFile` directive specifies a shared memory file that the `httpd` processes may use to track the state of the `JServ` processes.

Distribution of JServ Requests

The following process explains how `mod_jserv` selects the `JServ` engine to handle a request:

1. An `httpd` process is started.

2. `mod_jserv` creates a list of available JServs, with extra entries for JServs with a weight value greater than 1 (for example, `JServ2` in the example above, as specified by `ApJServBalance JServ_set JServ2 2`).
3. An `httpd` daemon receives a servlet request and hands it to `mod_jserv`.
4. `mod_jserv` selects the JServ engine that will handle the request.
 - a. `mod_jserv` checks to see if the request is part of a current session. If so, it uses the `ApJServRoute` directives to find the JServ that handled the other requests for that session.
 - b. If the request is not part of a session, `mod_jserv` selects an engine based on the process ID of the `httpd` process and the number of entries in the list of available JServs, as follows:

`JServ_id to handle the request = httpd_pid % (modulo) number of JServs in the list`

This method distributes requests across the available JServ engines fairly evenly.

Using Single Thread Model Servlets

Oracle recommends that you write your servlets to implement the `SingleThreadModel` (STM) interface. An application that was modified to implement the STM interface demonstrated a 25% improvement in **response time**.

It is also much easier to manage database connections with STM servlets. The database connection can be set up in the `init()` method of the servlet, and closed in the `destroy()` method. When executing the servlet's `doGet()` or `service()` method, you need not be concerned with obtaining a database connection. You can also manage database connections with JDBC connection caching.

There are three parameters in the `zone.properties` file that impact the performance of STM servlets in particular. These govern:

- The minimum number of servlet object instances that will be generated and available after the servlet class is loaded
- The maximum number that can be generated
- The number that should be generated if the available instances are insufficient

Because it is very costly to generate instances while the system is running, Oracle recommends that you set your minimum to equal your maximum value. The

optimum value depends somewhat on how many connections your database server can handle. This should be split among the JServ processes, as follows:

$$\text{Total DB connections} / \text{Number of JServ processes} = \text{Number of STM servlet instances per process}$$

See [Chapter 3, "Sizing and Configuration"](#) for suggestions on determining the right number of JServ processes for your application, and ["How to Perform Load Balancing"](#) on page 5-4 for the steps to configure them. Suppose you've determined that you want 10 servlet instances per process. The capacity settings in the `zone.properties` file would be:

```
singleThreadModelServlet.initialCapacity = 10
singleThreadModelServlet.incrementCapacity = 0
singleThreadModelServlet.maximumCapacity = 10
```

Warning: The value for `singleThreadModelServlet.maximumCapacity` in the zone properties file must be at least as large as the value for `security.maxConnections` in the `jserv.properties` file. If it is not, and the number of requests sent to the JServ process exceeds the maximum capacity, requests will fail.

What is OracleJSP?

OracleJSP 1.1.0.0 is Oracle's implementation of the Sun Microsystems JavaServer Pages 1.1 specification. Some of the additional features it includes are custom JavaBeans for accessing Oracle databases, SQL support, and extended data types. See the *Oracle9i Application Server Overview Guide* in the Oracle9i Application Server documentation library for detailed descriptions of the features.

Tuning OracleJSP Pages for Performance

This section explains how you can improve OracleJSP pages' performance.

Impact of Session Management

In general, sessions add performance overhead; they consume about 0.5 KB of resident memory. You must turn off sessions if you do not want a new session to be created with each request. By default, sessions are enabled in OracleJSP pages, so if

they are not being used, turn them off by including the following line at the top of the page:

```
<%@ page session="false" %>
```

If you are going to use sessions, ensure that you explicitly close them. If you don't, they will linger until they time out (the default value for session timeout is 30 minutes). To close a session manually, use the `session.invalidate()` method.

See the *OracleJSP Developer's Guide and Reference* in the Oracle9i Application Server documentation library for more information on configuring OracleJSP pages.

Developer Mode

Another parameter that has a significant effect on performance is developer mode. It is a useful feature for debugging during development, but it degrades performance. The default value is true, so you will need to set it to false in the `jserv.properties` file as follows:

```
servlet.oracle.jsp.JspServlet.initArgs=developer_mode=false
```

With developer mode set to true, OracleJSP and the servlet engine examines every request to determine whether to reload or retranslate the page or application. With developer mode off, only the first request is examined.

In a test using JDK 1.2 with 50 users, 128 MB heap, and the default TCP settings, the performance gains with developer mode off were 14% in **throughput**, and 28% in average **response time**.

Buffering

If an OracleJSP page is not using any features that do not require resetting the buffer (such as error pages, `contentType` settings, forwards, etc.), disabling the JSP page buffer will improve performance. This is because memory will not be used in creating the buffer, and the output can go directly to the browser. Use this page directive to disable buffering:

```
<%@ page buffer="none" %>
```

The default size of an OracleJSP page buffer is 8 KB.

OracleJSP Performance Tips

The configuration actions below can enhance the performance of your OracleJSP pages.

Caching database connections

Since the performance cost of creating database connections is high, it is more performant to use a cache of connections. If you use a cache of database connections, then the OracleJSP application can get a connection from the cache and return it when it is finished.

Configuring the statement batch value

The JDBC driver accumulates a number of execution requests (the batch value) and passes them to the database to be processed at the same time. You can configure the batch value to control how frequently processing occurs.

Caching JDBC statements

Cache executable statements that are repeatedly used, to avoid re-parsing, statement object re-creation, and recalculation of parameter size definitions.

Pre-fetching rows

During a query, pre-fetch multiple rows into the client to reduce round trips between the database and the server.

Caching rowsets from the database

Cache small sets of data that are accessed frequently and do not change often. This is not as beneficial for large data sets, since they consume more memory.

Invoking static includes

To invoke static includes, use the page directive:

```
<%@ include file="/jsp/filename.jsp" %>
```

Static include creates a copy of the file in the JSP, thereby affecting its page size. This is useful in avoiding trips to the request dispatcher (unlike dynamic includes, which must go through the request dispatcher each time). However, file sizes should be small to avoid exceeding the 64 KB limit of the service method of the generated page implementation class.

Invoking Dynamic Includes

To invoke dynamic includes, use the page directive

```
<jsp:include page="/jsp/filename.jsp" flush="true" />
```

This directive is analogous to a function call, and therefore does not increase the page size of the JSP. However, a dynamic include increases the processing overhead since it must go through the request dispatcher. Dynamic includes are useful for including other pages without increasing page size.

Index

A

Apache JServ Protocol 1.2, 5-2
ApJServBalance, 5-6
ApJServManual, 5-5
ApJServMount, 5-5
ApJServRoute, 5-6
ApJServShmFile, 5-6
architecture
 JServ, 5-2
 Oracle 9i Application Server, 1-11

C

caching
 database connections, 5-10
 SSL, 3-4
capacity, 1-7
concurrency
 defined, 1-2
 limiting, 1-8
concurrent executing users, 3-2
concurrent users, 3-2, 4-4
 MaxClients and, 3-2
connection caching, 5-7
contention, 1-5
CPU
 application requirements, 3-3
 average utilization, 3-3
 insufficient, 1-5
 statistics, 2-2, 2-3
 usage, 2-2
cron, 2-9

D

database connection, 5-7
demand rate, 1-6, 1-7
developer_mode, 5-9

E

ExtendedStatus, 2-7

F

functional demand, 1-7

G

graceful shutdown, 2-12

H

harsh shutdown, 2-12
hash
 defined, 1-2

J

JDBC, 5-7
JServ
 described, 5-2
 load balancing, 5-4
 process start-up time, 5-3
 processes, load balancing, 5-4
 starting and stopping processes, 5-4
 threads per, 3-8
JServ Protocol 1.2, 5-2

jserv.conf, 2-10
jserv.properties, 5-4
JSP, 5-8

K

kernel memory requirements, 3-5

L

latency
 defined, 1-2
 first-request, 5-3
 network, 3-2
load balancing, 5-4
load variances, 1-9
logging, 4-7

M

MaxClients
 concurrent users and, 3-2
 configuring, 4-6
 increasing, 2-8
memory usage, 3-5
mod_jserv, 5-2, 5-6
mod_oprocmgr, 5-4
mod_status, xi, 2-6
monitoring
 CPU usage, 2-2
 httpds processes, 2-6, 2-8
 JServ processes, 2-10
 server, 2-8
 server side status, 2-6
 server, automating, 2-9
mpstat, 2-2, 2-3
mpstat utility, 2-3

N

Network Monotor, 2-4

O

Oracle 9i Application Server 8i
 architecture, 1-11

oracle.jsp.jspServlet, 5-3

P

performance goals, 1-8, 3-1
protocol
 Apache JServ 1.2, 5-2
 HTTP/1.1, 4-8
 SSL, 3-4

R

repository, defined, 5-3
response time, 1-5
 defined, 1-2
 goal, 1-8
 improving, 1-3
 peak load, 1-9
 sizing and, 3-2

S

sar utility, 2-2
scalability
 defined, 1-2
 monitoring, 2-2
security.allowedAddresses, 5-5
security.maxConnections, 3-8
server statistics, 2-6
server-side status information, 2-6
server-status, 2-6
service time, 1-3
 defined, 1-2
servlet
 database connection and, 5-7
 engine, 5-2
 pre-loading classes, 5-3
 SingleThreadModel interface and, 5-7
 zone properties file, 5-3
servlets.startup, 5-3
sessions
 JServ processes and, 5-6
 SSL and, 4-7
SetHandler, 2-6
shutdown, 2-12

SNOOP, 2-4
SSL
 defined, 3-4
 performance cost, 3-4
 session caching, 4-7
statistics
 CPU, 2-2
 server, 2-6, 2-8
status reports, 2-6

T

think time
 defined, 1-2
 resources and, 3-2
thread
 limit, 3-8
 migrations to other processes, 2-4
throughput
 defined, 1-2
 demand limiter and, 1-7
 increasing, 1-4

U

unit consumption, 1-7
uptime, 2-6
users, concurrent, 3-2
utilities
 mpstat, 2-3
 sar, 2-2

W

wait time
 contention and, 1-5
 defined, 1-2
 parallel processing and, 1-4
 percentage of time spent, 2-4

Z

zone, defined, 5-3
zone.properties, 5-7

