# Load Balancing and Fault Tolerance in Distributed Web Servers

## Introduction

As web applications become more complex and handle increasing amounts of traffic, it is important to ensure scalability, high availability, and fault tolerance. Load balancing is a key strategy for distributing traffic among multiple servers to avoid overloading a single machine. Additionally, fault tolerance helps maintain availability in the event of a server failure.

In this document, we will explain how to set up load balancing and fault tolerance for a distributed web application, specifically a Task Manager servlet, using multiple instances of Apache Tomcat, Nginx as a load balancer, and a few optional features such as session persistence and database synchronization.

---

## 1. Setting Up Multiple Instances of Apache Tomcat

### 1.1 Installing Apache Tomcat

Apache Tomcat is an open-source web server that implements Java Servlet and JSP technologies. You can download the latest version from the official Apache Tomcat website. Once downloaded, you extract the archive and move the extracted files to a suitable location, such as `/usr/local/tomcat`.

### 1.2 Running Multiple Instances of Tomcat

To achieve load balancing, we will run two or more instances of Apache Tomcat on different ports. This requires creating separate copies of the Tomcat installation and configuring each instance to run on unique ports.

1. **Create multiple instances:** After installing Apache Tomcat, create two or more copies of the Tomcat installation folder (e.g., `tomcat1`, `tomcat2`).
2. **Configure ports:** For each Tomcat instance, you must modify the `server.xml` configuration file found in the `conf` directory. Update the HTTP connector port for each instance to a unique port (e.g., `8080` for the first instance, `8081` for the second).

### 1.3 Deploying the Task Manager Application

Once the Tomcat instances are configured, deploy the Task Manager servlet in each instance. This involves placing the `TaskManager.war` file (or extracted files) in the `webapps` folder of each Tomcat instance. After deployment, you can start each Tomcat server and access the application on different ports.

---

# 2. Installing and Configuring Nginx as a Load Balancer

## 2.1 Installing Nginx

Nginx is a widely used open-source web server that can act as a reverse proxy and load balancer. To install Nginx, download the package from the official website or use your operating system's package manager.

## 2.2 Configuring Load Balancing in Nginx

Once Nginx is installed, you need to configure it to balance the load between the multiple Tomcat instances. This is done by editing the `nginx.conf` file, typically located in the `/etc/nginx` directory. The key settings for load balancing include defining the upstream servers (your Tomcat instances) and setting up the load balancing rules.

For example, you define the upstream servers as your Tomcat instances running on different ports (e.g., `localhost:8080` and `localhost:8081`). You can choose different load balancing strategies such as round-robin (default), least connections, or IP hash.

## 2.3 Session Persistence (Sticky Sessions)

To ensure that requests from the same client are always routed to the same Tomcat instance, you can enable session persistence. Nginx provides the `ip_hash` directive, which routes requests based on the client's IP address. This ensures that users remain on the same Tomcat server for the duration of their session.

## 2.4 Testing Load Balancer Setup

Once the load balancer is configured, start Nginx and Tomcat. You can now test the load balancer by accessing your web application via Nginx and observing how requests are distributed across the Tomcat instances. Simulate multiple client requests using browser tabs or a tool like `curl`.

---

# 3. Implementing Fault Tolerance

### 3.1 Handling Server Failures

Fault tolerance ensures that when one Tomcat instance fails, the load balancer redirects traffic to another available instance. In Nginx, this can be achieved using health checks. Nginx will periodically check the status of each Tomcat server, and if a server goes down, it will be automatically removed from the pool of available servers.

You can simulate server failures by stopping one of the Tomcat instances and observing how Nginx reroutes traffic to the remaining active instances.

### 3.2 Automatic Failover

In addition to handling sudden server failures, Nginx can automatically reroute traffic when a server becomes unavailable. You can configure the load balancer to redirect requests to a secondary server if the primary server fails to respond within a specified timeout period.

---

# 4. Database Synchronization (Optional)

If your Task Manager application stores persistent data (such as task lists), it is important to ensure database consistency across the different Tomcat instances. This is typically done by configuring a single central database that all Tomcat servers can access. Alternatively, you can use database replication to ensure that each instance has a synchronized copy of the database.

### 4.1 Centralized Database

In a centralized database setup, all Tomcat instances connect to the same database server. This ensures that data remains consistent no matter which instance handles a request. Popular databases like MySQL, PostgreSQL, and MongoDB support centralized setups.

### 4.2 Database Replication

For more complex setups, you can implement database replication, where each Tomcat instance maintains a local copy of the database that is kept in sync with the others. Databases like MySQL provide built-in replication features. However, this approach can be more complex to set up and manage.

---

# 5. Optional Features

### 5.1 User Authentication

You can enhance the Task Manager application by adding user authentication. Using Java Servlets, implement login and logout functionality with session management. When a user logs in, store their session ID in a cookie. Each request should include the session ID, and the application should validate it to ensure the user is authenticated.

### 5.2 REST API

A REST API allows for easy integration with other services and provides a more modern interface for the Task Manager application. You can implement REST endpoints using `doGet`, `doPost`, `doPut`, and `doDelete` methods in your servlets. The Task Manager can expose endpoints such as `/tasks` to retrieve all tasks, `/tasks/{id}` to retrieve a specific task, and `/tasks` (POST) to create new tasks.

### 5.3 AJAX for Front-End Updates

To improve user experience, you can use AJAX in the front-end to dynamically update the task list without requiring a page reload. This can be implemented with JavaScript by sending asynchronous requests to the server and updating the DOM with the response data.

### 5.4 HTTPS for Secure Communication

Securing communication between clients and the server is critical for production applications. You can enable HTTPS in Nginx by obtaining an SSL certificate (via Let's Encrypt, for example) and configuring the Nginx server to serve content over HTTPS. This ensures that data is encrypted in transit and protected from interception.

---

# 6. Testing and Evaluation

After setting up the load balancer and deploying the Task Manager application on multiple Tomcat instances, thoroughly test the system's performance, fault tolerance, and user session handling.

### 6.1 Load Testing

Simulate multiple users accessing the application to evaluate how effectively Nginx balances the load between Tomcat instances. Monitor each server's resource usage and response times to identify any bottlenecks.

## 6.2 Fault Tolerance Testing

Deliberately take one Tomcat instance offline and observe how Nginx handles the failure. Ensure that users are automatically rerouted to the remaining instances without experiencing downtime.

## 6.3 Database Synchronization

If database replication is implemented, test data consistency by submitting tasks via different Tomcat instances and verifying that all instances have access to the same data.

---

# Conclusion

This document outlines the steps required to set up load balancing and fault tolerance for a distributed web application using Apache Tomcat and Nginx. By running multiple instances of Tomcat and using Nginx to distribute traffic, you ensure that the system can handle increased loads and continue to operate in the event of server failures. Optional features such as user authentication, REST APIs, and database replication further enhance the application's robustness.