# Performance Testing Workload Modeling: Ramp-Up Duration and User Load Calculation (Technical Guide)

## 1. Overview of Workload Modeling in Performance Testing

Workload modeling is a crucial step in performance testing that defines how users interact with a system over time. The key components include:

- **User Load (Concurrency):** The number of virtual users (VUs) interacting with the system simultaneously.

- **Ramp-Up Strategy:** The rate at which users are introduced into the system.

- **Steady-State Duration:** The time for which the system is tested at peak load.

- **Ramp-Down Strategy:** The rate at which users exit the system.

These components ensure that the test accurately simulates real-world user behavior and helps in identifying system bottlenecks.

---

## 2. Factors Influencing Ramp-Up Duration & User Load Calculation

### A. Business and SLA Requirements

Workload modeling must align with business goals and Service Level Agreements (SLAs). Key business requirements include:

1. **Expected Peak Load:** How many concurrent users does the system need to handle?

2. **Transactions Per Second (TPS):** Business-defined throughput requirements.

3. **Response Time SLA:** Acceptable response time for critical transactions.

**Example:**

For an **e-commerce platform** expecting **10,000 concurrent users** on **Black Friday** with an SLA of **≤ 2s response time**, the performance test should simulate:

- Peak load of 10,000 users.

- Average **500-1000 transactions per second (TPS)**.

- A response time **≤ 2 seconds**.

---

### B. Real-World Traffic Analysis

To simulate realistic user load, we analyze historical traffic data from:

1. **Monitoring Tools:**

- **Google Analytics, New Relic, Dynatrace, AppDynamics** provide insights into:
  - Peak traffic patterns.
  - Session durations.
  - User abandonment rates.

2. **Web Server & API Logs Analysis:**
   - Extract peak requests from web server logs:

   grep "POST /checkout" access.log | wc -l

   - Use log analysis tools like **Splunk or ELK Stack** to visualize trends.

3. **User Distribution Patterns:**
   - Traffic varies across time zones and business hours:
     - **Enterprise applications:** 9 AM - 6 PM peak.
     - **E-commerce:** Traffic surges in the evening and weekends.
     - **Streaming services:** Peak usage at night.

---

**C. System Capacity & Scaling Considerations**

- **Fixed Infrastructure:** If no auto-scaling is enabled, a slow ramp-up prevents overloading.

- **Cloud Auto-Scaling (AWS ASG, Kubernetes HPA):**
  - The ramp-up should match the **auto-scaling threshold** (e.g., every 5 minutes).
  - Example: If AWS **Auto Scaling Group (ASG)** scales every **2 minutes**, the ramp-up should introduce **100 users every 2 minutes**.

**Example: AWS ASG Scaling Thresholds**

| Metric | Scaling Trigger |
|---|---|
| CPU > 75% | Scale up by 2 instances |
| TPS > 1000 | Scale up by 1 instance |
| Memory > 80% | Scale up by 3 instances |

---

**3. How to Calculate Ramp-Up Duration?**

The ramp-up duration is determined based on:

- **Total Users (N):** Number of virtual users.

- **Ramp-Up Time (T):** Time to reach full load.

- **Ramp-Up Rate:** The speed at which users are added.

## A. Linear Ramp-Up Formula

Ramp-Up Rate=Total Users / Ramp-Up Time (T)

Example:

- **5,000 users over 20 minutes**: 5000\20=250 users per minute

  - **JMeter Configuration:** Set **Thread Group Ramp-Up = 1200 seconds (20 min).**

---

## B. Stepwise Ramp-Up (Batch Injection)

Rather than adding users continuously, they are injected in **batches** to prevent early failures.

## Example: Stepwise Ramp-Up (Batch Injection)

- **1,000 users every 5 minutes** until full load is reached.

| Time Interval | Users Added | Total Users |
|---|---|---|
| 0 - 5 min | 1,000 | 1,000 |
| 5 - 10 min | 1,000 | 2,000 |
| 10 - 15 min | 1,000 | 3,000 |
| 15 - 20 min | 1,000 | 4,000 |
| 20 - 25 min | 1,000 | 5,000 |

## JMeter Configuration:

- Use **"Ultimate Thread Group"** to set **ramp-up batches**.

---

## C. Arrival Rate-Based Ramp-Up (Little's Law)

If the system is **TPS-driven**, we calculate concurrency using:

Concurrent Users=TPS×Response Time

## Example Calculation

- **Target TPS: 1,000 transactions per second**.

- **Average Response Time: 2 seconds**.

- **Required Concurrent Users:** 1000×2=2000 concurrent users

- **JMeter Configuration:**

o   Use **"Constant Throughput Timer"** to control the **TPS rate**.

---

## 4. How to Determine Total User Load?

There are multiple approaches:

### A. Production Traffic Analysis

- Extract peak **Concurrent Users** from **New Relic, Dynatrace, AppDynamics**.

- Web Server Logs:

  cat access.log | grep "GET /home" | wc -l

### B. Capacity Planning & SLA

- **Max Load Capacity**: If the database supports **2,000 TPS**, and avg response time is **3 sec**, then: 2000×3=6000 concurrent users

- **SLA-Based Load Testing**: Adjust user load to find the maximum sustainable number of users under the SLA constraints.

---

## 5. Performance Test Workload Models

| Model | Description | Use Case |
|-------|-------------|----------|
| **Peak Load** | Simulates the highest expected user load. | SLA validation |
| **Step Load** | Gradually increases load in steps to find breaking points. | Capacity planning |
| **Spike Load** | Sudden large user injection. | Flash sales, DDoS testing |
| **Soak Test** | Sustains load for long periods. | Memory leaks, stability testing |

---

## 6. Example Workload Scenario

**E-Commerce Black Friday Load Test**

- **Expected Peak Users: 50,000**

- **TPS Estimate: 4,000 TPS**

- **Avg Response Time: 3 sec**

- **Concurrent Users Calculation:** 4,000×3=12,000

**Ramp-Up Strategy:**

- **First 10% (5,000 users) in 5 min**

- **Next 30% (15,000 users) in 15 min**

- **Final 60% (30,000 users) in 30 min**

**Steady State: 1 hour Ramp-Down:** Gradual drop over **15 minutes**.

---

**7. Conclusion**

- **Linear and stepwise ramp-up** are most commonly used.

- **Arrival rate-based ramp-up** aligns with real-world TPS expectations.

- **JMeter can implement these strategies using Thread Groups & Timers**.

---

There are **multiple models** for calculating **ramp-up duration** in performance testing. While **linear, stepwise, and arrival rate-based ramp-ups** are common, **several other methods exist** based on real-world usage patterns, system architecture, and scalability constraints.

Below are **other 3 ramp-up models**, each with detailed calculations and examples.

---

**4. Spike Ramp-Up (Sudden Traffic Burst)**

- **All users are injected within a short time** (e.g., 1-2 minutes).

- Used for **stress testing, DDoS simulation, and flash sales scenarios**.

**Example Calculation:**

- **5000 users injected in 2 minutes**: 5000\2=2500 users per minute

| Time Interval | Users Added | Total Users |
|---------------|-------------|-------------|
| 0 - 1 min | 2500 | 2500 |
| 1 - 2 min | 2500 | 5000 |

**JMeter Configuration:**

- Use **Stepping Thread Group** to inject all users quickly.

**Use Case:**

- **Ticket booking, Black Friday sales, breaking news websites**.

---

**5. Wave-Based Ramp-Up (Oscillating Load)**

- Load fluctuates **up and down in waves** to simulate **real-world traffic fluctuations**.

- Useful for **handling scaling challenges in cloud environments**.

**Example Calculation (Waves Every 5 Minutes)**

| Time Interval | Users Added | Total Users |
|---------------|-------------|-------------|
| 0 - 5 min | +1000 | 1000 |
| 5 - 10 min | -500 | 500 |
| 10 - 15 min | +1500 | 2000 |
| 15 - 20 min | -1000 | 1000 |
| 20 - 25 min | +2500 | 3500 |

**JMeter Configuration:**

- Use **Custom Thread Groups** to define wave-like traffic.

- Ideal for **IoT applications, streaming platforms, and mobile apps**.

---

<u>S a n t h o s h   K u m a r   J</u>

**6. Ramp-Up with Auto-Scaling Considerations**

- Designed for **AWS Auto Scaling, Kubernetes Horizontal Pod Autoscaler (HPA), Azure VM Scale Sets**.

- Users are added in **sync with scaling triggers**.

**Example Calculation (Scaling Every 5 Minutes)**

| Time Interval | Users Added | Total Users | Scaling Action |
|---|---|---|---|
| 0 - 5 min | 1000 | 1000 | 1 instance added |
| 5 - 10 min | 1500 | 2500 | 2 instances added |
| 10 - 15 min | 2500 | 5000 | 4 instances added |
| 15 - 20 min | 5000 | 10,000 | 8 instances added |

**JMeter Configuration:**

- Set **gradual ramp-up matching cloud scaling intervals**.

**Use Case:**

- **Cloud-native microservices that rely on Kubernetes auto-scaling**.

---

**Comparison of Ramp-Up Models**

| Model | Best For | Ramp-Up Strategy |
|---|---|---|
| **Linear Ramp-Up** | General applications | Fixed rate increase |
| **Stepwise Ramp-Up** | Legacy systems | Users injected in batches |
| **Arrival Rate-Based Ramp-Up** | APIs, backend services | TPS-driven load |
| **Spike Ramp-Up** | Flash sales, stress testing | Sudden load increase |
| **Wave-Based Ramp-Up** | Streaming, IoT systems | Oscillating traffic patterns |
| **Auto-Scaling Aware Ramp-Up** | Cloud-based microservices | Matches scaling policies |

---

**Example Workload Scenario: E-Commerce Black Friday Sale**

A **retail website** expecting **50,000 concurrent users** during **Black Friday** should use a **hybrid ramp-up strategy**:

1. **Stepwise Ramp-Up** for first **10,000 users** (to warm up caches, CDNs).

2. **Spike Load** (20,000 users in 5 min) simulating a flash sale.

3. **Wave-Based Load** (±5,000 users every 10 min) to mimic real-world purchase behavior.

**Test Configuration**

| Phase | Users | Ramp-Up Time | Load Model |
|---|---|---|---|
| **Warm-Up** | 10,000 | 10 min | Stepwise |

| | | | |
|---|---|---|---|
| **Flash Sale** | 20,000 | 5 min | Spike Load |
| **Steady-State** | 50,000 | 1 hour | Linear |
| **Load Drop** | -20,000 | 15 min | Wave-Based |

**Conclusion**

1. **Not all performance tests use simple linear ramp-up.**

2. **Stepwise and wave-based ramp-up models** better replicate **real-world traffic.**

3. **Cloud-based tests** should align with **auto-scaling policies.**

4. **Arrival rate-based models (TPS-driven) are best for APIs.**

5. **Spike ramp-up is essential for stress testing and event-driven traffic surges.**