

## Gangigunta, Amarnath (CORP)

---

**From:** Gangigunta, Amarnath (CORP)  
**Sent:** Sunday, May 29, 2022 6:02 PM  
**To:** Gangigunta, Amarnath (CORP)  
**Subject:** Performance Engineering Basics

### Performance testing Vs Engineering :

- Performance Testing is a quality check of the application in terms of the application's response handling capacity. Performance Testing verifies how a system will perform under production conditions and anticipate issues that might arise during heavy load conditions. On the other hand, Performance Engineering aims to design the application by keeping the performance metrics in mind and also to detect and resolve issues early in the development cycle.

### Software architecture:

Before going to explain the types of architecture firstly you need to understand the different layers and foundation of software architecture. There are four types of layer:

- **Presentation Layer:** This layer is responsible to display the user interface and manage user input.
- **Application Layer:** Application layer (also known as Business Layer) has all the business logic and policies. The application layer is a bridge between the presentation layer and the data layer.
- **Data Layer:** This layer is responsible for storing the data.
- **Service Layer:** This layer is responsible to define and implement the service interface and logic. Service layer communicates with the application layer.

**One Tier:** MP3 player, MS Office

**Two tier:** Client(Presentation and app) – server(Database) application

### Performance metrics:

**Client Tier:**

- TCP connection time
- HTML resources load time
- CSS files load time
- Images load time
- JavaScript file load time
- HTTP response time
- HTTP response status

**Web tier:**

- Cache Hit length
- Request Queued
- Number of HTTP connections

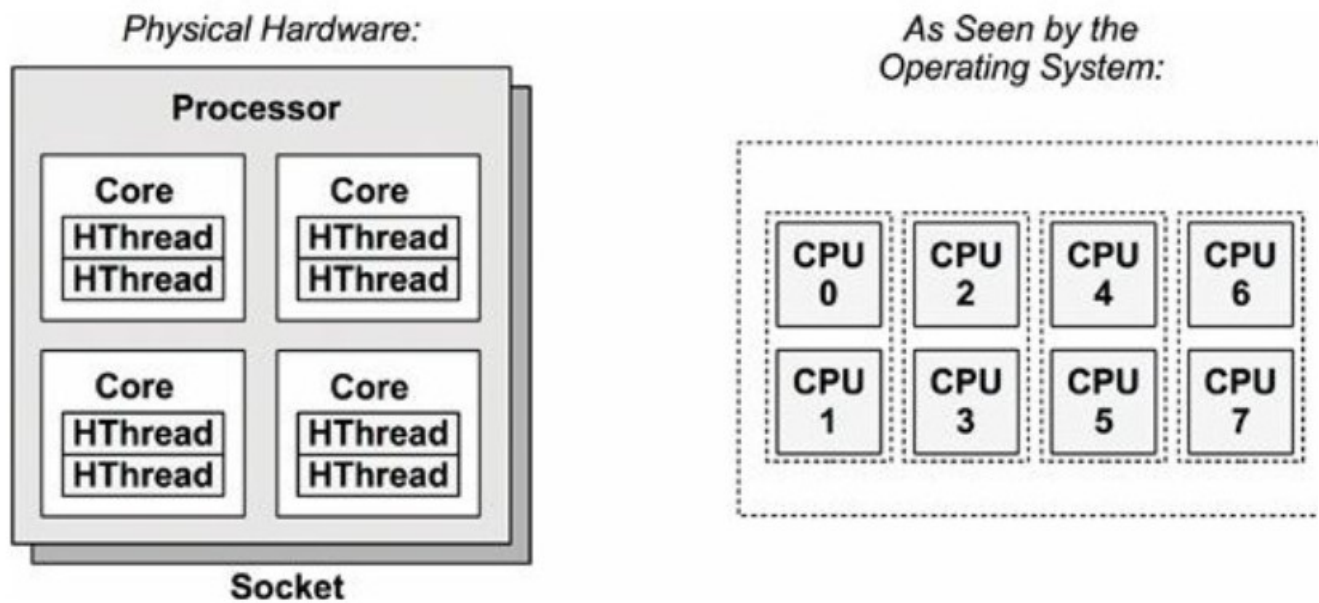
**Application tier:**

- Connection Time
- Connection Wait Time
- Total Threads
- Memory Use
- Active Transactions
- Transactions suspended
- Rolled Back transactions
- Timeouts
- Servlet Errors

**Database Tier:**

- Query throughput
- Query execution performance
- Queue Length
- Connections
- Buffer pool usage

**Physical Vs Logical CPU:**



Now that we know the required terms, a **Physical CPU** means the actual Physical Core is present on a processor. The multithreading being enabled on the core or not doesn't change the count of Physical CPUs present on a Processor. **In short, Physical CPU correlates actual Physical Cores present on a processor. In the above diagram there are 4 Physical Cores present on the Single Processor.**

**Logical CPU** internally refers to the ability of each core doing 2 or more tasks simultaneously. This is achieved by enabling hyperthreading on the cores. Each single physical core is divided into multiple logical cores by enabling hyperthreading on them. In the above diagram, although there are only 4 Physical cores, the system is tricked to look at it as 8 Logical CPUs, by enabling hyperthreading on each core.

The above diagram correlates to below details :

1 Processor, 4 Physical cores, 8 Logical Cores. In other words - 1P4C - with 2 threads per core.

**Agent Vs Agent Less monitoring tools:**

## Difference Between Agent-based and Agent-less Performance Monitoring Tool:

- **Proprietary Software:**

- *Agent-based:* Need to install on each server which is under monitoring
- *Agent-less:* No need to install on each server

- **Central System:**

- *Agent-based:* Along with agents, there is one central system (machine) named Master o This machine communicates with all the agents and collects the data.
- *Agent-less:* There is only one central system which directly communicates with the server (under monitoring) and collects the data.

- **Installation:**

- *Agent-based:* It requires the installation of a central system with the additional installation on each server.
- *Agent-less:* It requires the installation of a central system only.

### Factors that affect performance:

#### Code:

- memory leaks
- array bound errors
- inefficient buffering
- too many processing cycles
- a larger number of HTTP transactions
- too many file transfers between memory and disk
- inefficient session state management
- thread contention due to maximum concurrent user
- poor architecture sizing for peak load
- inefficient SQL statements
- lack of proper indexing on the database tables
- an inappropriate configuration of the servers

#### Network:



- Older or unoptimized network infrastructure
- Slow web site connections lead to network traffic and hence poor response time
- Imbalanced load on servers affecting the performance

#### **Garbage Collection:**

A Garbage Collector is a Java program which tracked the referenced (live) objects and allowed them to stay in the heap memory whereas the memory of the unreferenced (dead) objects is reclaimed and reused for new object allocation. This method of reclaiming the unused memory is known as Garbage Collection. In C++, the programmer is responsible for managing memory, while in Java, the Garbage Collector handles it.

Garbage collection is related to memory management and helps to improve the performance of a system. Garbage collection in a system should not be too quick or too late. Too many GC cycles degrade the performance of the system and causing the spike in CPU whereas a delay in GC cycles leads to memory leaks.

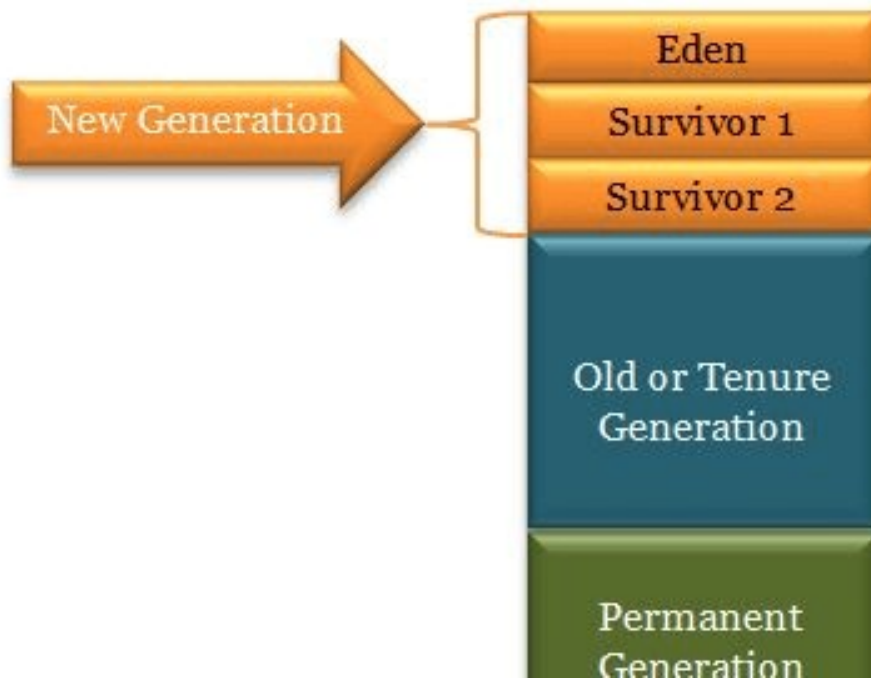
#### **Some important terms:**

1. *Live Object*: The object which is referenced by another object
2. *Dead Object*: The object which is not referenced by any other object and unreachable
3. *Daemon Thread*: Garbage collection process is carried out by Daemon Thread
4. *System.gc()*: It is used to invoke the garbage collector and on invocation, the garbage collector reclaims the unused memory space. The `System.gc()` is a static method.

#### **Memory heap:**

##### **What is Memory Heap?**

Heap is a portion of memory which is used for dynamic allocation. The blocks of the memory in the heap are frequently allocated and freed based on the status of the object. The allocation and release of memory in the heap take place in an arbitrary order and hence the unreferenced or dead objects exist in the heap. The memory of the unreferenced objects is reclaimed by the Java program called Garbage Collector.



#### Garbage Collection Process:

##### **Promotion:**

There is one more process which takes place along with all above-mentioned activities i.e. Object Promotion. After each minor GC, the objects become old and their age keeps on increasing. When aged objects reach a certain age threshold they are promoted from the young generation to the old generation. This process is called Promotion.

In our example you can see object A, G, J and O are promoted to the old generation after 'n' number of GC executions. The promotion process leads to fill old generation memory space. One important point to note here is that some of the objects like String and Array objects are directly created in the tenured generation space.

#### Types of GC:

**Serial:** Uses only one thread for GC

**Parallel:** multiple threads for GC

**Concurrent mark sweep (CMS) GC:**

The Concurrent Mark Sweep (CMS) garbage collector collects the tenured generation. It is also called concurrent low pause collector. It attempts to minimize the pauses due to garbage collection by doing the garbage collection work concurrently with the application threads. Due to this reason, the CMS uses more CPU than other GCs. If you can allocate more CPU for better performance then CMS garbage collector is a preferred choice over the parallel collector. CMS garbage collector uses multiple threads to scan the

**G1 (Garbage First) Garbage Collector** : G1 GC provides benefit where multi-processor machines along with large memory space are available.

**Stop the world event:**

**Stop the World Event** - All minor garbage collections are "Stop the World" events. This means that all application threads are stopped until the operation completes. Minor garbage collections are *always* Stop the World events.

**The Old Generation** is used to store long surviving objects. Typically, a threshold is set for young generation objects and when that age is met, the object gets moved to the old generation. Eventually the old generation needs to be collected. This event is called a *major garbage collection*.

**Heap Dump:**

JVM heap dump is an informative data (sometimes it is called a snap) of Java heap memory which contains low-level details about java objects and classes allocated in the Java heap at the moment when the dump is fetched or snapshot is triggered. Generally, a heap dump is triggered after the full GC ran so that the dump contains the information about the remaining objects in the heap. A heap dump contains the information

**Thread Dump:**

Similar to Heap Dump, Thread Dump is a snapshot of the status of all the threads at a particular time

## 1. Thread:

A thread is a series of statements which can be processed in parallel in the same program to achieve concurrency. All the Java threads are managed by JVM (Java Virtual Machine) which are mapped with operating system thread, also called as Native Thread. Threads help to allow multiple activities within a process and act like lightweight processes. Each thread has a unique identifier, name, and category.

## 2. Multithreading:

Java is a multi-threaded application that allows multiple threads to execute at any particular time. A web application uses tens to hundreds of threads to process a large number of concurrent users and achieve multithreading. The problem occurs when two or more threads utilize the same resources and leads to a deadlock situation.



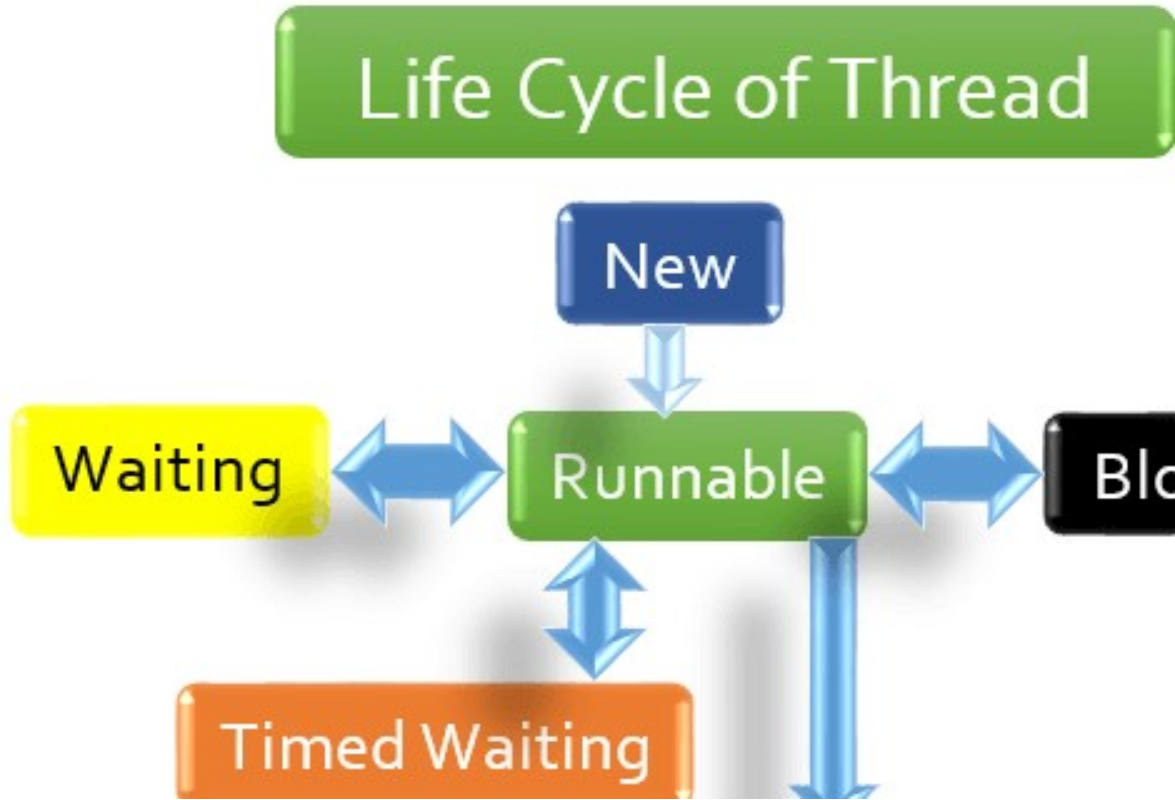
### 3. Deadlock:

Deadlock is a situation when two or more threads are waiting for the other threads to complete order to complete their own tasks. To identify the root cause of the deadlock, a thread dump is:

### 4. Thread Dump:

Re-iterating the definition again. Thread Dump is a snapshot of the status of all the threads at helps to find out what every thread in the JVM is doing at a particular point in time. Are some deadlock? Are all the threads running? Which threads are waiting for the blocked resources et questions can be answered by analysing in the thread dump.

#### Life Cycle of a Thread:



Thread Dump Analyzer: **fastThread**



## Blocked Thread Scenario

To continue the thread dump analysis, now we will move to the next part which will cover the analysis of blocked threads. The meaning of a blocked thread is that a thread locks a resource for a long time and restricts other threads from using that particular resource, hence the restricted threads are moved into the Blocked state as Blocked Threads. Let's try to understand with an example. Refer to the below fastThread dashboard.

## Deadlock Scenario

Deadlock is a situation when the locked resources are required by some other threads who also require resources which are needed by other threads. Such a cyclic condition restricts threads from completing their execution, resulting in none of the threads completing its execution and the system is halted. Refer to the below fastThread dashboard.

### **AWR :**

Automatic Workload Repository (also called AWR) report provides information about the database bottleneck.

### **Dump vs AWR:**

The difference between the dump (heap or thread) and Automatic Workload Repository report is that dump reports are the actual snapshots taken at a particular time whereas AWR report is a comparison of snapshots taken at different timestamps.

## AWR Features:

The Automatic Workload Repository report provides:

- Wait-events causing a delay
- Highlights query taking long Elapsed time or execution time
- Report on CPU utilization and memory usage
- Blocking sessions and many other important stats which we will be discussed in later posts

### **Server Tuning Tips:**

**Web server:** Apache, ASP, IIS

