

# HP Global Workload Manager—Improving server CPU utilization technical overview



Introduction .....	2
HP Virtual Server Environment (VSE) overview .....	3
VSE components .....	4
Partner integration .....	5
WLM 3.0 and HP gWLM 1.1 comparisons .....	5
HP gWLM overview .....	5
HP gWLM customer benefits .....	6
Improving resource utilization through automation .....	6
Easily managing resource policies across multiple servers .....	7
Centralized monitoring and reporting .....	7
Customer use cases .....	7
Consolidating environments and stacking applications .....	7
Improving utilization of existing capacity .....	7
Improving application performance .....	8
Providing overdraft protection for mission-critical workloads .....	8
Freeing CPUs for additional workloads .....	8
Utilizing resources on Serviceguard clusters .....	8
Adjusting resource estimates based on what is actually used .....	8
HP gWLM features and functions .....	9
For more information .....	10

## Introduction

Today, most servers are highly underutilized. Traditional IT environments have typically been configured as silos where resources are aligned around an application or business function. Capacity is fixed, resources are over-provisioned to meet peak demand, and systems are complex and difficult to change. ERP, CRM, and web-based applications, for example, typically operate in silos. Many reports estimate average server utilization at approximately 30% in UNIX® environments. Costs are based on owning and operating the entire vertical infrastructure, even when underutilized. Customers often have one application per server, and they size that server for peak loads of typically three to five times the average utilization.

This strategy is inefficient and often results in:

- Fixed capacity and cost not aligned with business needs
- Under-utilized and over-provisioned server resources
- A complex architecture that is difficult to change

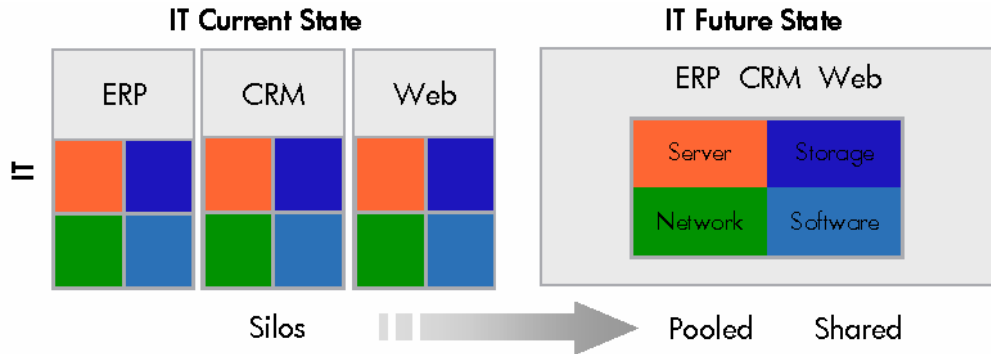
Virtualization is the abstraction of hardware through software. It enables a business to pool and share IT resources, as shown in the following graphic, so utilization is optimized and supply automatically meets demand. By pooling the IT resources that form your infrastructure and sharing these across applications and business processes, virtualization enables a rapid response to change. Virtualization is a cornerstone of HP's approach to helping customers realize the promise of becoming an Adaptive Enterprise—where business and IT are synchronized to capitalize on change.

Server virtualization enables administrators to optimize the usage and simplify the management of single- or multiple-server environments by configuring them as reusable pools of resources. Server virtualization offers many benefits; however, its effectiveness is maximized when the element server virtualization functions are combined with a powerful method for allocating and optimizing resources across and between virtual servers. The HP Virtual Server Environment (VSE), an integrated solution for both HP Integrity and HP 9000 server platforms, allows organizations to achieve a greater return on their IT investments by optimizing server resource utilization on a real-time basis according to business priorities. Within a VSE, virtual servers automatically grow and shrink based on the service-level objectives set for each application they host. Through tight integration with high availability, partitioning, and utility pricing, the VSE allows organizations to maintain service levels in the event of unexpected downtime and to pay for spare capacity on an as-needed basis.

A virtual server environment is more agile and responsive to change, enabling customers to:

- Double their resource utilization
- Maintain continuous server levels
- Pay only for what they use

Figure 1.



For HP-UX environments, the HP Virtual Server Environment is built around the only automated, goal-based policy engines available for UNIX—HP-UX Workload Manager (HP-UX WLM) or HP Global Workload Manager (gWLM)—which allow for consolidation of multiple applications on a single server without compromising performance. For example, HP-UX WLM or HP gWLM can automatically move CPU resources from one virtual partition to another, or HP-UX WLM can automatically activate and deactivate utility pricing CPUs to optimize the cost of IT resource utilization. These resource shifts can occur purely based on CPU utilization or time of the day—the majority of HP customers use that as an indicator for their application performance. Optionally, HP offers other service-level objectives, such as number of active users or response time thresholds, that can be determined by an administrator. In addition, in a clustered, high-availability environment, HP-UX WLM can react to the service level objective of an application package that has failed over from another partition or server and allocate the necessary resources on-the-fly to ensure that service levels are maintained.

HP gWLM is a key component of the HP Virtual Server Environment, helping organizations pool and share IT resources to improve utilization and align supply with demand. Let's take a closer look at VSE and workload management.

## HP Virtual Server Environment (VSE) overview

HP VSE delivers the highest degree of integrated virtualization to ensure customers the speed and agility to adapt. HP VSE optimizes resource utilization by consolidating and virtualizing server resources, where each server intelligently adapts (grows or shrinks in size) based on business priorities and service-level objectives (SLOs).

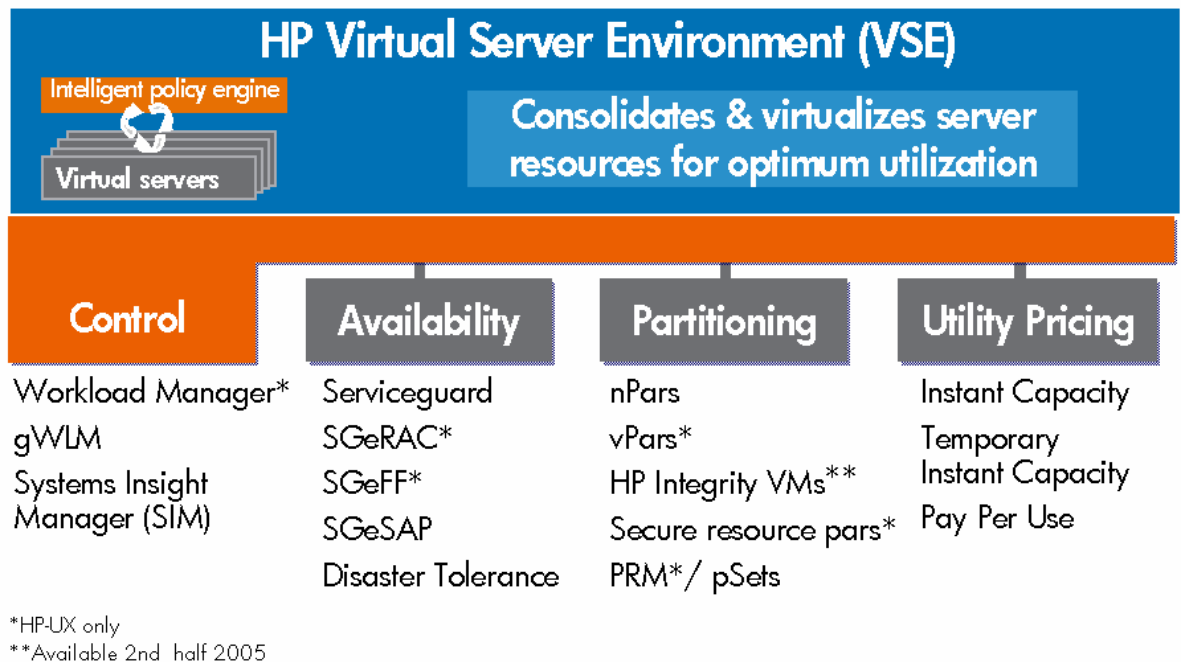
HP offers a broad spectrum of virtualization solutions that allow customers to choose the most appropriate path and optimization focus of their IT infrastructure resources. HP defines these strategies in the following terms:

- Element Virtualization—*Optimize utilization of individual servers, storage, and other resources to meet demand within a single application environment or business process* (Partitioning, Clustering, Instant capacity, Pay per Use, Rapid Deployment, Workload Management, HP Storage Solutions—Enterprise Virtual Array (EVA)).
- Integrated Virtualization—*Optimize multiple infrastructure environments to automatically meet service-level agreements* (Virtual Server Environment (VSE), Storage Grid, HP BladeSystem, Consolidated Client Infrastructure (CCI), Disaster Tolerant Solutions (DTS))
- Complete IT Utility—*Optimize all heterogeneous resources so supply meets business demand in real time* (Grid Computing, HP Managed Services (Strategic Outsourcing))

## VSE components

The following figure shows the components of VSE.

Figure 2.



nPartitions (nPars) are hardware-based, electrically isolated partitions that enable customers to deploy multiple operating systems on a single server.

Virtual partitions (vPars) are separate operating system instances on the same nPartition or system, with OS, application, and resource isolation. HP-UX 11i Virtual Partitions enable you to dynamically move CPU power between vPars as your workload requirements change. vPars offers single CPU granularity.

Resource partitions, including Secure Resource Partitions, PRM groups, and processor sets, allocate resources to specific applications and users within an operating system. They offer fully dynamic application CPU or sub-CPU and percent memory granularity.

## Partner integration

HP VSE has many significant partner-integration capabilities. VSE is application-transparent—applications can take advantage of the benefits of VSE with no need for application changes. The most common configuration is to change resource allocation based on the workload's CPU utilization. In addition, HP-UX Workload Manager has ISV toolkits available for Oracle, BEA, SAS, and Apache. These toolkits extract application specific information (such as a number of active users for Oracle or queue length for BEA) and use that information to make better resource allocation decisions. HP Serviceguard is certified to work with thousands of applications so that every application can take advantage of high availability.

HP has developed several HP Virtual Server Environment Reference Architectures for HP-UX 11i that include all the components required for a BEA WebLogic or Oracle Database environment, including nPars, vPars, PRM, Psets, HP-UX Workload Manager, Serviceguard (including SGeRAC) and Instant Capacity. The VSE Reference Architecture information includes cookbooks, configuration files, scripts, testing results and management best practices. These pre-tested configurations allow our customers up to 60% faster implementation of these environments.

The HP Virtual Server Environment works with HP Systems Insight Manager (HP SIM). HP SIM offers additional third-party integration and integration with higher level management products such as HP OpenView Operations.

## WLM 3.0 and HP gWLM 1.1 comparisons

HP-UX Workload Manager (WLM) is the optimal solution for a line of business (LOB) consolidation in which the LOB owns its servers but relies on IT to manage them. It allows a company to run applications on a modest number of systems with usually fewer than 20 partitions. WLM gives companies the ability to prioritize applications and create measurable application SLOs. Today, HP-UX Workload Manager also offers a deeper integration with HP Serviceguard and Utility Pricing offerings. HP-UX Workload Manager is optimized to work on a single system.

HP gWLM is optimal for a central IT consolidation that owns and manages servers for other businesses. It allows many applications to be run on a large number of systems, with many partitions on those systems. It provides an IT utility to the LOB, allowing IT to rent resources to the LOB, based on application sizes and own/borrow/lend, fixed-allocation, and CPU-utilization models. For more details on these models, refer to the "Customer use cases" section. HP gWLM requires a central management server that is integrated with HP Systems Insight Manager. This central management server allows HP gWLM to work in a multi-OS environment by managing both HP-UX and Linux servers running the HP gWLM agent software.

## HP gWLM overview

To form the Virtual Server Environment, gWLM integrates with virtualization offerings such as resource management groups, partitioning, instant capacity, and clustering. HP gWLM allows you to monitor workloads in different partitions and set policies to adjust the amount of resources a workload is allocated based on which applications are busy or idle, according to the rules you set to govern the sharing of resources.

HP gWLM allows many applications to run on a virtualized environment consisting of multiple systems and partitions. Within this environment, resources can be moved to the applications as needed. HP gWLM monitors the resource requests of the environment and reallocates resources automatically based on business policy. If business conditions change, gWLM can automatically reassign resources to handle the new conditions, optimizing resource utilization while maintaining continuous service levels. gWLM does this by automatically resizing virtual partitions and resource partitions.

gWLM continuously monitors this environment to determine how resources should be allocated. This information is consolidated and stored on the central management server and provides the data for various types of automated reports. In addition to the many reports it offers, HP gWLM also has a real-time monitoring facility that shows customers actual resources, requested resources, actual utilization, target utilization, minimum size, and maximum size of partitions for multiple systems.

## HP gWLM customer benefits

HP gWLM is the policy engine that monitors real-time resource requirements and automatically reallocates resources to the virtualized servers. Dynamic resource allocation in a multi-OS environment can double resource utilization. Simple, easy-to-use, policy management means continuous service levels are maintained while improving the utilization.

HP gWLM benefits customers in the following ways:

- It improves resource utilization through automation.
- Centralized management and HP Systems Insight Manager (SIM) interface make it easy to use.
- Centralized monitoring and reporting show which applications are using resources and how often.

## Improving resource utilization through automation

The past strategy for complex IT structures has been to run applications on separate, physical servers, and hope that these servers can handle the resource demand. This siloed infrastructure limits resource use to whatever is available on any particular server and restricts capacity movement. Over-provisioning is often needed to achieve performance goals, and siloed infrastructures are sized for peak workload, leading to IT resources that are significantly underutilized. For many companies, overall utilization of IT resources can be less than 30%, while some applications are still not able to meet performance requirements.

Rather than sizing siloed infrastructures for peak workload, HP gWLM takes these siloed applications and manages the resource management once they are consolidated onto virtualized servers. The virtualized servers are sized for normal workload with a pool of shared resources to handle the peaks across multiple applications. In a virtualized environment, IT capacity can be automatically moved and excess capacity can be shared among applications based on real-time demand. One virtual server can, in essence, borrow unused resources from another virtual server, leading to more fully utilized IT resources, and lower expenses for software licenses, support, space, power, installation, and integration. HP gWLM automatically reallocates resources to workloads based on:

- Own/borrow/lend policy—This type of policy allows you to set the amount of CPU resources a workload owns, the minimum amount of CPU resources a workload should ever have, and the maximum amount of CPU resources a workload should ever have. Servers with such a policy are guaranteed the owned amount of CPU when needed. The minimum CPU allocation allows you to specify how much a workload can lend, if owned resources are not needed, and the maximum allows you to specify how much the workload can borrow, if the amount owned is insufficient. If a workload has lent out CPU resources and the workload becomes busy, the workload re-acquires those lent-out CPUs immediately.
- Fixed-allocation policy—This type of policy guarantees that a workload has a fixed (constant) amount of CPU resources.
- CPU utilization policy—This type of policy has a target based on utilization. With a CPU utilization policy, gWLM attempts to keep a workload's CPU utilization within a target range by adding CPU resources when the workload is using too much of its current CPU allocation and reducing resources when utilization drops below the target utilization range.
- Custom Policy (Goal-based Management)—This type of policy allows you to provide your own metric. gWLM then manages an associated workload, adjusting CPU allocation as needed, based on how the value of its metric compares to a specified target.

## Easily managing resource policies across multiple servers

In a traditional data center environment, each silo of infrastructure is customized, so specific knowledge is required to implement changes. Because a virtualized infrastructure is standardized with a logical interface, changes can be implemented more easily and quickly (often dynamically).

Using gWLM, customers can easily implement and manage an automated virtualized environment across many servers within a datacenter. HP gWLM can be managed in an integrated fashion from HP Systems Insight Manager (SIM)—the central point of administration for complete resource life-cycle management for multi OS environments—through the Central Management Server (CMS). The CMS consists of a web-based UI and a repository for storing the configuration and performance data. gWLM provides the ability to manage different types of virtualized servers in a unified manner. This single product manages both resource allocation within an operating system image (full and sub processor) and can allocate resources between multiple partitions on the same server.

## Centralized monitoring and reporting

The key control point for HP VSE is the intelligent policy engine, HP gWLM, which performs real-time monitoring and assessment of the resource utilization, and then advises and acts in accordance with the pre-set policies. By monitoring application performance or CPU utilization and shifting resources from idle workloads to busy workloads, HP gWLM enables you to achieve the same or better performance levels with less infrastructure. Pre-enabled, real-time and historical reports show that workloads received resources when they were required. Because a virtualized infrastructure is shared, gWLMs monitoring and reporting functionality enables internal charge-back. The pre-enabled reports also show which application instances borrow or lend excessive resources and show that an application instance got its guaranteed resources whenever it required them.

Virtualization, and the automation of a virtualized datacenter, is a relatively new technology in the datacenter. Because many customers may not be comfortable with the automation, HP gWLM offers an advisory mode in addition to its active management mode that can be configured to send an alert to HP Systems Insight Manager instead of automatically moving resources and/or provide a report on gWLM's advice. This allows customers to view actual resources, requested resources, actual utilization, target utilization, and an approximation of the actions gWLM would have performed.

## Customer use cases

### Consolidating environments and stacking applications

When combined with partitioning technologies, gWLM makes it easy to create flexible, consolidated environments when consolidating in these three ways:

- Consolidating multiple production environments on the same server using partitions
- Consolidating test/dev and production environments on the same server using partitions
- Consolidation through application stacking within the same OS image

Here are several examples of how HP gWLM can be used in a customer environment.

### Improving utilization of existing capacity

If the workloads on your vPars are growing, eventually, demand will outgrow existing resources and you will be forced to purchase additional CPUs or migrate to new hardware to handle the higher demand. By using gWLM to automatically migrate CPUs to vPars where they are needed most, you can better use existing capacity, delaying or eliminating the need for additional hardware.

## Improving application performance

By migrating CPUs from idle vPars to busy ones, each workload has more computing power available to handle peak loads without purchasing additional resources.

## Providing overdraft protection for mission-critical workloads

In a static environment, vPars are sized to handle the projected peak load. But if demand unexpectedly exceeds forecasts, insufficient CPU resources can cause unacceptable performance problems. With gWLM, if a vPar gets an unexpected surge in demand, that vPar can borrow CPUs from another vPar, according to the policies you put in place. Because most workloads' periods of peak demand are infrequent, there is a high likelihood that CPUs will be available.

## Freeing CPUs for additional workloads

By sharing resources, gWLM allows you to achieve the same performance with fewer CPUs. If you already have a large vPar installation, when you turn on sharing with gWLM, you can then use some of the CPUs for additional workloads.

## Utilizing resources on Serviceguard clusters

In a situation where two servers are being used with a Serviceguard package—one server running applications and the other idle for use as a hot standby—the resources of the standby server are completely unused except in the rare event of the failure of the active server. If the standby server is a vPar, other vPars running low-priority applications on the same physical platform can use the resources owned by the standby server. In the event of a failure, the resources are automatically migrated from low-priority applications back to the standby instance that is now performing work.

## Adjusting resource estimates based on what is actually used

In a static environment with one application per server that is not on a shared platform, planning estimates are heavily padded. When estimating what is needed to run important applications, it is typical to err on the side of caution. However, on a shared platform with gWLM managing resources, if estimates are incorrect, gWLM allocate resources to your workloads based on what the workloads are actually using. A pool of resources can be shared among a large number of applications. Assuming that all of these applications will not be at peak load at the same time, resources float between applications as they are needed. Allocating resources based on resource demand you are actually experiencing enables you to use shared headroom, reducing the amount of padding needed.



# HP gWLM features and functions

HP gWLM features provide functions that enable customers to reallocate resources for their applications automatically, based on centralized management policies and see the results of this action through centralized monitored reporting capability.

---

Automation of resource allocation	<p>gWLM manages vPars, pSets, and Fair Share Schedule (FSS) groups on HP-UX—11i v1 for HP 9000 server and 11i v2 for HP Integrity and HP 9000 servers.</p> <p>gWLM manages CPU affinity for SUSE Linux Enterprise Server 9 and Red Hat Enterprise Linux AS 4 on HP Integrity servers.</p>
Easy-to-use centralized management policies	<p>A centralized management server makes management convenient. A web-based GUI integrates with HP Systems Insight Manager, providing easy-to-use central policy management and centralized monitoring and real-time reporting, as well as historical utilization, demand, and resource allocation data.</p> <p>It provides out-of-the-box support for common, easy-to-understand configurations—own/borrow/lend policy, fixed allocation policy, and CPU utilization policy. Discovery of virtual and resource partitions is automatic, and a single policy can be applied to multiple resource partitions through the central management server. Most sites only need several variations of factory-supplied policies, and customers can begin managing a vPars environment immediately after installing software. Customer test results showed customers could actively manage their environments within 15 minutes.</p> <p>Access to gWLM configuration and real-time and historical reports is role-based. System administrators have full access, and anyone with read-only access can only view reports.</p>
Centralized monitoring and reporting	<p>Advisory mode and reports such as real-time reports and historical reports, ensure that gWLM is configured properly. Application owners can see results of policies and automation of resource allocation. CPU resources can be guaranteed according to fixed policy, own/borrow (PolicyMin, owned), and utilization (PolicyMin). Audit reports prove each workload got what was guaranteed and show when workloads got more resources than the guaranteed minimum.</p>

---

## For more information

For more information, visit <http://www.hp.com/go/gwlm> or contact your local authorized HP reseller or HP sales office.

© 2005 Hewlett-Packard Development Company, L.P. The information contained herein is subject to change without notice. The only warranties for HP products and services are set forth in the express warranty statements accompanying such products and services. Nothing herein should be construed as constituting an additional warranty. HP shall not be liable for technical or editorial errors or omissions contained herein.

Linux is a U.S. registered trademark of Linus Torvalds. UNIX is a registered trademark of The Open Group.

5983-0505EN, 03/2005

