



***Make IT Service Monitoring Simple and Proactive
with Intelligent Thresholding and Alerting***

White Paper

Restricted Rights Legend

The information contained in this document is confidential and subject to change without notice. No part of this document may be reproduced or disclosed to others without the prior permission of eG Innovations Inc. eG Innovations, Inc. makes no warranty of any kind with regard to the software and documentation, including, but not limited to, the implied warranties of merchantability and fitness for a particular purpose.

Trademarks

Microsoft Windows, Windows 2000, Windows 2003, Windows 2008, Windows Vista, Windows XP and Windows 7 are either registered trademarks or trademarks of Microsoft Corporation in United States and/or other countries. Citrix®, MetaFrame® XenApp® and ICA® are registered trademarks of Citrix Systems, Inc. in the US and other countries.

The names of actual companies and products mentioned herein may be the trademarks of their respective owners.

Copyright

© Copyright eG Innovations. Specifications subject to change without notice.

Introduction

IT managers often complain about two main types of problems with monitoring and management tools:

- Firstly, after they install the software, they start to receive many “false” alerts. A false alert refers to a situation in which the monitoring tool indicates a problem but the IT manager determines that there is no real problem in the network. Thousands of alerts can result distract IT administrators, prevent them from focusing on the real issues that can impact IT service quality.
- Secondly, to avoid false alerts, IT managers have to define threshold values for the different metrics collected by the monitoring tool. A threshold is a limit set in the monitoring tool for the metric, so that if a metric crosses this value, an alert is raised. In a large enterprise, a monitoring tool that provides visibility into the different network, server, and application tiers can collect millions of metrics. Having to set thresholds manually for each and every metric is a very time consuming, monotonous exercise. As a result, enterprises end up spending a lot of time and money having consultants tune thresholds manually.

Hence, what administrators need from a monitoring and management system is the ability to make thresholds simple to configure and accurate to enforce, so there are few false alerts. The following sections define how the eG Enterprise IT service monitoring solution from eG Innovations addresses this key requirement of IT managers using a combination of automated thresholding and intelligent alerting. By doing this without requiring a lot of manual intervention, eG Enterprise makes it simple to implement IT service monitoring in an enterprise, and yet deliver proactive alerts that are essential for ensuring that the service level expectations of users are met.

Defining Static Thresholds for Metrics

As indicated earlier, thresholds are upper and lower bounds that determine whether a metric is performing to expectation or not. Every time the actual value of the metric falls outside the prescribed limits, the monitoring system detects an abnormality.

Depending on the metric being collected, upper bounds are appropriate for some metrics, while lower bounds are appropriate for others. For example, a lower bound or minimum threshold is applicable for the free disk space metric. If the value drops below the lower bound, an alert will be generated. In contrast, an upper bound or maximum threshold is applicable for the CPU usage metric of a server. If the value exceeds the upper bound, an alert will be generated. In some cases, both lower and upper bounds may be appropriate. For instance, if the number of users accessing a server is much higher or much lower than normal, it could be an indicator of a problem.

For many metrics, thresholds can be set statically. For instance, based on the service level expectations and agreements, IT managers can set thresholds for metrics such as network availability and latency. Application availability and response time can also be handled in the same manner. For example, availability should be 100% whenever the metric is measured. If not, a violation should be detected. Likewise, a network latency of several seconds is usually an indicator of a problem, no matter what time of day the measurement is made at. Figure 1 shows the CPU usage metric for a server compared with its static threshold value.

Reducing False Positives in Network Monitoring

“The million-dollar question is: How can you reduce false positives for counters that tend to fluctuate a lot? I have seen some administrators try to reduce the sampling frequency in an effort to reduce false positives. Indeed, this technique may reduce false positives, but it still has the same result. Counters that fluctuate a lot will still produce false positive alerts.”

—by Brien Posey, SearchNetworking.com

<http://searchnetworking.techtarget.com/tip/Reducing-false-positives-in-network-monitoring>

Stop Monitoring Tools from Crying Wolf

“In my opinion, a poorly tuned monitoring server is as bad or worse than no server monitoring at all. At least with no monitoring you are less likely to become complacent. If you don't have a car alarm and live in a bad neighborhood, you'll probably be more careful to put away valuables and lock your doors. But if you have a car alarm that goes off every time another car drives by, you will naturally start to ignore it over time.”

—By Kyle Rankin, Data Center Systems Management,
<http://searchdatacenter.techtarget.com/tip/Stop-server-monitoring-tools-from-crying-wolf>

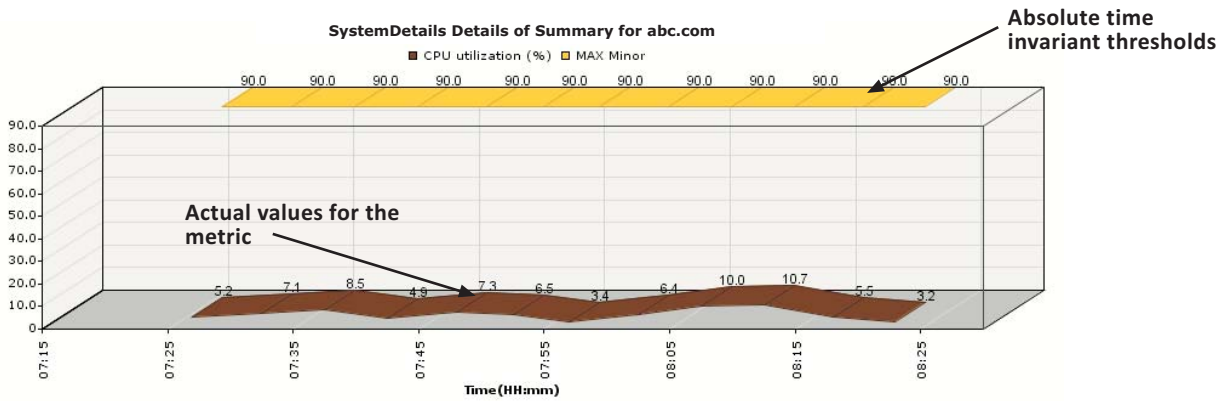


Figure 1: A graph indicating the value of a measure and its absolute threshold

Thresholds can also be set based on industry standard best practices. For example, a rule of thumb when tuning an Oracle database server is that the database dictionary cache hit ratio should be 90% or more. If the hit ratio falls below this value, it indicates a need to tune the database server. This is another example where a threshold is set statically, without considering the time of day when the measurement is being made. eG Enterprise includes pre-specified threshold values for many metrics based on industry standard best practices.

Often, there is a need to set different threshold levels to map to different levels of severity of problems. For example, when the space usage of a disk is close to 90%, an IT manager would like to receive a minor alert. When the metric's value crosses 95%, the IT manager would like to receive a major alert and when the value crosses 99%, the IT manager would like to receive a critical alert. To support such requirements, eG Enterprise allows administrators to set different threshold levels for the same metric (see Figure 2). Based on the value that is crossed, an alert with the appropriate priority is generated by eG Enterprise. Multiple levels of threshold settings allow proactive alarms to be generated when a metric is slightly out of conformance, and a severe alarm to be generated when the problem worsens.

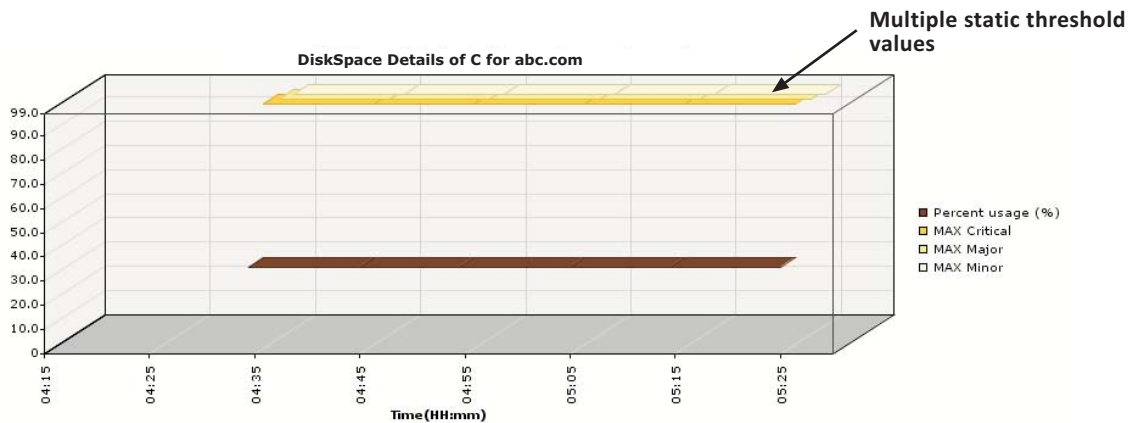


Figure 2: Multiple thresholds set for the "Percent usage" metric for disk space of a Windows server

Defining Automatic, Self-Adjusting Thresholds for Metrics

In infrastructures where a metric varies with time, a static absolute threshold value cannot serve as a reliable basis for judging performance. For example, consider a web server hosting a web site. The number of TCP connections to the web site (i.e. the current connections measure of the TCP test in figure 3) could be rather high on a particular day and low on another. Similarly, it could be high during the working hours and low during the nights. In such situations where measurement values change with the time of the day, it is very difficult to set accurate maximum and minimum limits

The Growing Adoption of Behavior Learning Tools

"By 2015, 10% of organizations will implement behavior learning tools to supplement existing IT operations management tools, up from fewer than 1% today."

—By Debra Curtis and David Williams
 "Seeking Patterns With IT Operations Management Tools",
 Gartner Report, 2010.

Behavior-learning Tools are Making Inroads in IT Operations Management

"These tools enable IT organizations to proactively identify and isolate faults and performance issues in the IT infrastructure by consolidating fault and performance data from a wide range of sources, establishing a "normal" behavior pattern or profile, which is then automatically analyzed and compared with newly gathered data to detect subtle anomalies in real time."

—By David Williams

"Tools Alone Won't Enable Proactive Management of IT Operations",
Gartner Report, 2010.

manually. In such cases, the threshold value for this metric also has to be time variant.

Even when a metric is not time variant, its value may change from one server to another. For example, a high-end datacenter server may be able to handle hundreds of users, whereas a low-end standard server may be able to handle only a few tens of users. In such cases too, it is extremely laborious and time consuming to determine what the normal values are for each and every server.

To handle such situations, eG Enterprise includes an auto-thresholding capability. Using past history of the values of the metric, eG Enterprise uses tried and tested statistical quality control techniques to analyze

past values of the metrics and to automatically set the upper and lower bounds (i.e., baselines) for each of the metrics. The thresholds are set based on the eG Enterprise system learning the behavior of the underlying infrastructure automatically. For example, the threshold values for a metric between 9am-10am tomorrow are based on the value of the metric for the same time period over the past days (the number of days to be looked at in the past is configurable).

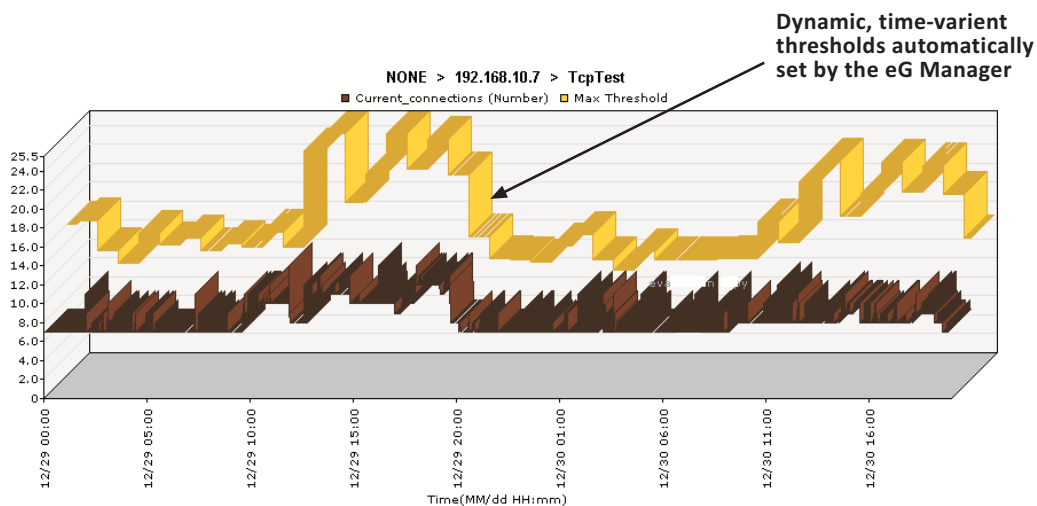


Figure 3: A performance graph illustrating how automatic threshold values are computed by the eG manager based on the actual values of a metric

Figure 3 shows an example. In this case, the number of TCP connections to a server is tracked over time. The yellow line in this figure denotes the time-varying threshold that the eG manager has automatically computed for this metric. The number of days for which the behavior of the infrastructure is observed to compute the baselines automatically is configurable, with two weeks being the default.

As is evident from Figure 3, eG's auto-thresholding capability ensures that like the metric value, the threshold also is time varying. Whenever a deviation from the baseline (upper or lower thresholds) is detected, an alert is triggered. Since the baseline is set automatically, using this technique ensures that IT managers are informed of problems well before they become critical enough to impact the end user experience.

Recommendations for Behavior Learning in IT Operations Management

- ✓ IT operations must associate normal and abnormal behaviors with either good or bad behaviors, and the value of the results increases when aligned to a business service model.
- ✓ Behavior patterns need time to be established; therefore, the value diminishes if pattern-based tools are introduced into IT environments that change erratically with unplanned, unmanaged changes.
- ✓ Good-quality results require good-quality data. If the underlying, supporting data sources are sending poor data, then the behavior patterns will be created with bad data.

—By Debra Curtis and David Williams

"Seeking Patterns With IT Operations Management Tools",
Gartner Report, 2010.

Automatic thresholding is ideal for time varying metrics such as number of requests to a web server, the workload on a database server, queue lengths of requests waiting for processing, etc.

Even when thresholds are set automatically, an IT manager may want to choose a leniency factor for the thresholds. For example, an IT manager may want to allow for a 10% deviation from the norm. To accommodate such requests, eG Enterprise allows administrators to set a "sensitivity slider" for automatic thresholds. This slider should be specified as a multiple of the threshold value computed using statistical quality control (sqc). For example, consider the case of the "Free memory" measure, which is an indicator of the amount of free memory available on a server. Assume that on one of the managed servers, the free memory is known to decrease consistently and then grow back up (e.g., the operating system frees memory periodically). In such a scenario, the free memory threshold will be violated often (since the value decreases consistently), and this will result in a number of false alerts. In such a situation, the eG administrator can set the threshold to be a multiple of sqc - for example, if the minimum threshold is set to $0.7 * sqc$, it implies that the administrator has introduced a 30% leniency. That is, alerts are generated only if the free memory is 30% lower than what is the normal value. This capability allows administrators to fine-tune eG's relative thresholding capability to suit their specific requirements.

Multiple levels of thresholds are also supported when setting automatic thresholds – for example, a minor alert can be generated when a 10% leniency limit is crossed, a major alert generated when a 30% limit is crossed and a critical alert generated when a 50% limit is crossed.

Auto-Static Combination Thresholds

Automatic thresholds are ideal for metrics that are time variant. Often, the same metric may vary significantly from one server to another and from time to time. Consider a staging environment with a Citrix server. Typically, there is no load on the Citrix server and the automatic threshold is set accordingly. When someone logs in, the threshold will be breached and an alert may be raised by the system. This is a false alert because one user logging in does not signify a situation of interest to an IT manager. This scenario shows that while automatic thresholding reduces the effort involved in configuring the monitoring tool (because IT managers do not have to configure thresholds for every metric and server), it does not eliminate false alerts.

Therefore, eG Enterprise allows IT managers to use a combination of static and automatic thresholds. A static threshold applied along with an automatic threshold provides a realistic boundary that has to be crossed before an alert is to be triggered. An IT manager can now configure an absolute maximum and an automatic maximum threshold for a metric. eG Enterprise compares the actual measurement value with the higher of the two maximum thresholds, and generates an alert only when the higher threshold is violated. In the example of the staging Citrix server, the IT manager can set a static limit of say 10 sessions. Once this is done, only if the actual load exceeds 10 current sessions will an alert be generated, even if the auto-computed threshold is less than 10. If the auto-computed threshold is greater than 10, this value is used as the actual threshold.

Automatic Thresholds can be Incorrect as Well

"A Microsoft SCOM alert indicates that the "Terminal Services Active Sessions" metric is above the baseline. Active Sessions metric is above the calculated baseline. Current value is 1.33333333333333. How can I correct it and what does it mean?"

A posting on Microsoft System Center Forums, June 2009.

The Problem of Relying Exclusively on Automatic Thresholds

"In response to "Why do I get an alert when the terminal services active sessions metric exceeds 1.333333?":

There isn't a fix as such as it is behaving the way it should. It is just that self tuning thresholds are a pain in the backside and should be avoided, especially when low values are returned as a small value change can represent a large percentage change. They are great in theory but just don't work very well in practice.

Any time a customer is not happy with the results of a self-tuning threshold monitor – they should simply create a static threshold monitor. This is very basic and provides the best solution."

–by Graham Davies in Microsoft System Center Forum, 2009.

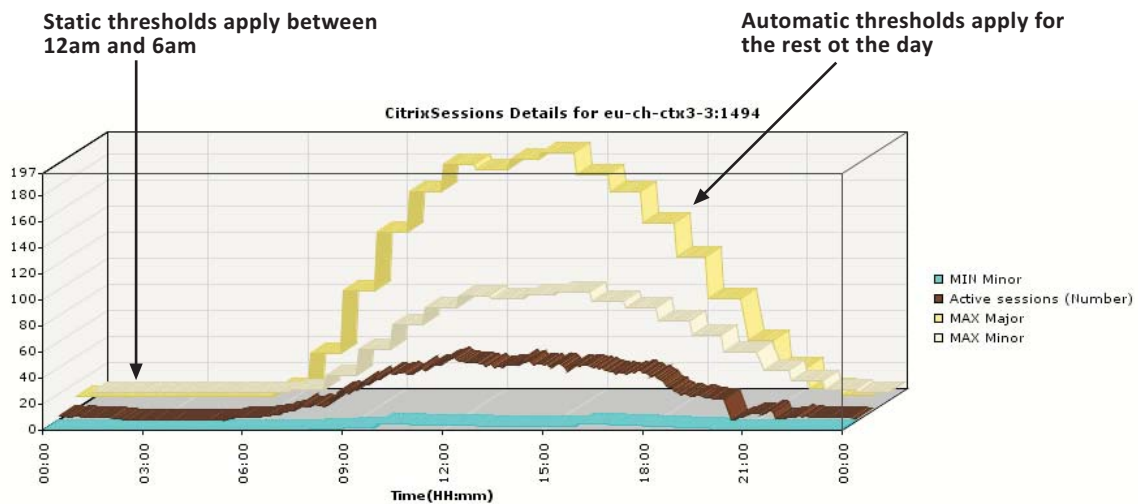


Figure 4: A performance graph showing number of Citrix sessions. An auto-static combination threshold is applied to this metric. In the morning hours, a static threshold is applied because the automatic threshold is lower. The static value ensures that alerts are not generated as long as the number of sessions stays below 5. During the day (6am onwards), the automatic threshold takes over.

Figure 4 shows an example of auto-static thresholds computed by the eG manager. In this example, the number of user sessions to a Citrix server is the metric under consideration. The brown line in the figure denotes the metric's value over time. The yellow lines represent the upper threshold values. Notice that from 12am to 6am, the threshold is static – with the minor value at 5 sessions and the major value at 10 sessions. Since the automatically computed value is less than both thresholds, the statically set threshold values apply in this case. From 6am onwards, as user activity increases, the automatically computed threshold value exceeds the static threshold value and hence, the auto-computed threshold applies.

As in the case with the maximum thresholds, if a static minimum and an automatic minimum threshold are specified, then eG Enterprise will generate alarms only when the current value falls below the lower of the two threshold settings.

The key benefits of eG's thresholding approach are:

- ✓ IT managers have the flexibility to choose different threshold policies for each and every metric. Different policies can be set for different types of servers and applications, and even for individual servers and applications.
- ✓ Threshold configuration is completely automated. IT managers do not have to sit and configure thresholds for every metric and application/server.
- ✓ There is no need to continuously tune thresholds as the IT infrastructure evolves, servers are added, applications are resized, etc. The monitoring system automatically learns and adapts to these changes.

Flexible Alarm Policies

Threshold configurations help determine when the state of a metric changes, but a threshold violation might not necessarily indicate a problem condition worthy of being reported to help desk. In other words, a single threshold violation might not always be reason enough for an alarm to be generated by eG Enterprise.

While a threshold policy determines how the thresholds for a metric are computed, an alarm policy determines when alarms are to be generated to inform IT managers about a problem. Depending on their criticality, different metrics may require different alarm policies. For instance, an instantaneous surge in the CPU usage of a system is a natural phenomenon in a production system. On the other hand, even a sporadic unavailability of a critical network router is a critical event that needs to be informed to the administrator. Alarm policies must also take into account the frequency of threshold violations of a metric. E.g., while an instantaneous surge of the CPU usage is not a cause for concern, a prolonged set of surges of the same metric may indicate a problem situation that must be corrected.

To accommodate different types of metrics, the eG Enterprise offers IT managers complete flexibility in setting alarm policies. IT managers can set individual alarm policies for each server, or each server group, or per server type (e.g., web server, database, application server, etc). An alarm policy is a combination of two parameters, **window size** and **number of crossings**. The **window size** represents the number of measurement values that are considered in determining the current state of a measurement. A crossing indicates a measurement being in violation of its threshold (i.e., a measurement value being lower than its lower threshold value, or a measurement value being higher than its upper threshold value). The **number of crossings** denotes the number of times a measurement has crossed its threshold.

Each of eG's alarm policies defines a window size and number of crossings (see Figure 5 below). For example, an immediate policy has a window size of 1 and number of crossings value of 1. This means only one measurement value is considered in determining the state of a measurement. If the current value exceeds the upper threshold limit, the measurement is said to be in an abnormal state, since number of crossings is 1. As its name indicates, this intelligent thresholding policy is ideal for cases where the administrator needs to be alerted immediately when an anomaly occurs. Metrics such as network / application availability can be monitored using this policy. For some other metrics, an administrator may not wish to be bothered about a sporadic threshold violation and may prefer to be alerted if a problem remains for a period of time. The standard alarm policy can be ideal for this, as it has a window size of 6, with number of crossings as 4.

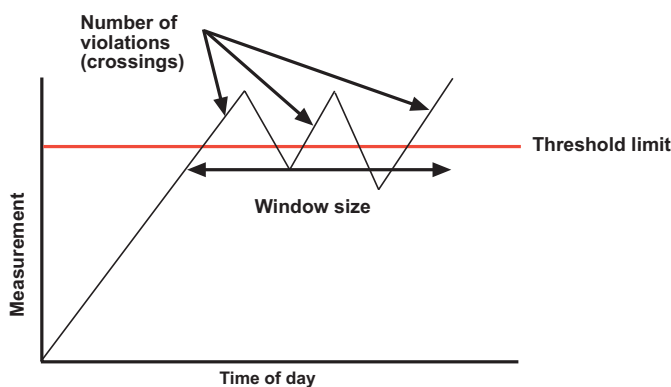


Figure 5: Defining Alarm Policies – This figure shows the concepts of “window size” and “number of crossings”

Choosing the Right Metrics is a Key to Being Proactive

Intelligent thresholding and flexible alarm policies are necessary, but not sufficient for a monitoring system to be effective. The metrics collected by the monitoring system are extremely important as well. If a monitoring system only collects availability and response time metrics, the metrics are not great early warning indicators of problems. This is because response time has an exponential distribution with load (see Figure 6) - i.e., as load increases, initially response time stays low, but as the load increases beyond the acceptable limit, response time shoots up dramatically even with a small variation in the load. This means that monitoring systems that use thresholds for response times are often not good at forecasting when problems are likely to occur.

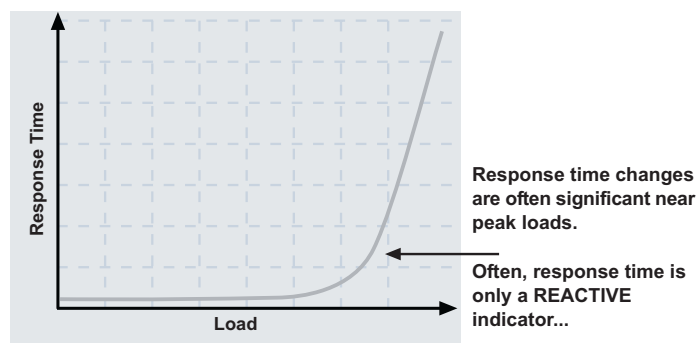


Figure 6: The variation of response time with load. The exponential relationship of response time with load means that response time is often not a good proactive indicator of problems.



A monitoring system that has domain knowledge built into it is often more effective than one that does not. For example, for a monitoring system monitoring VMware vSphere, it is essential that the monitoring system be able to look at CPU ready times of the virtual machines (VMs). By doing so, the monitoring system can determine times when the servers do not have sufficient CPU processing power. Early warning indicators can be provided to alert administrators as this issue starts to occur often. Likewise, tracking the number of requests waiting for I/O on a server can indicate a disk bottleneck, which if left unattended over time can result in catastrophic service outages.

Summary

In summary, to be effective, a monitoring system should include capabilities for thresholding intelligently, without requiring manual intervention, provide flexible policies for alerting administrators on alerts and be able to track key metrics in the underlying infrastructure. Having all three of these capabilities is key to enabling enterprises manage their IT services effectively. Early warning indicators provided by such systems can direct administrators to potential problems which if not fixed can have catastrophic consequences. The benefits of such a system are manifold. By being proactive, the monitoring system ensures provides administrators with the indicators they need to fix problems without impacting the business services they are responsible for supporting. Intelligent thresholding makes the configuration and implementation simple, thereby ensuring that the monitoring system can be up and running quickly in a cost-effective manner. On-going usage and maintenance of the monitoring system is also simplified by the intelligence built into the monitoring system. Since they receive fewer false alerts, administrators can focus their attention on the key problems in the infrastructure, rather than being distracted by a large number of meaningless alerts.

The eG Enterprise performance monitoring, diagnosis, and reporting solution from eG Innovations incorporates all of these key monitoring, thresholding, and alerting capabilities. For more information on eG Enterprise, please visit <http://www.eginnovations.com>.

For More Information

-  Read a whitepaper titled "[The eG Approach to Root-Cause Analysis](#)" to see how intelligent root-cause diagnosis can further simplify problem detection, troubleshooting and resolution in your IT infrastructure.
-  Listen to this webcast titled "[Managing N-Tiers without Tears](#)". This webcast discusses why managing multi-tier application infrastructures is a challenge and presents a detailed discussion on the key capabilities of eG Enterprise that make it an ideal solution for multi-tier IT infrastructures.

About eG Innovations

eG Innovations, Inc. (<http://www.eginnovations.com>) is a global provider of performance monitoring and triage solutions for virtual, physical and cloud-based IT infrastructures. The company's patented technologies provide proactive monitoring of every layer of every tier in the infrastructure, thereby enabling rapid diagnosis and recovery in enterprise and service provider networks. By ensuring high availability and optimum performance of mission-critical business services, eG Innovations' solutions help enhance customers' competitive positioning, lower operational costs and optimize the performance of their infrastructures.

USA

eG Innovations, Inc.

33 Wood Ave. South, Suite 600
Iselin, NJ 08830
Ph: (866) 526 6700

SINGAPORE

eG Innovations Pte Ltd

33A Tanjong Pagar Road
Singapore 088456
Ph : (65) 6423 0928
Fax : (65) 6423 1744

UK

eG Innovations UK Ltd.

3rd Floor,
126-134 Baker Street,
London W1U 6UE
Ph: +44 (0)20 7935 6721
+44 (0)12 7650 1590

Rest of Europe

eG Innovations, Europe

Montval
29 rue Leonard de Vinci
59700 MARCQ-EN-BAROEUL
France
Ph: +33 (0)3 66 64 06 16

INDIA

eG Innovations Pvt Ltd

2, Murali Street, Mahalingapuram
Chennai 600 034
India
Ph : (91) 44 2817 2801
Fax : (91) 44 2817 9041

Email : sales@eginnovations.com

Web : www.eginnovations.com

