

Table of Contents

<u>Service Level Management: Best Practices White Paper</u>	1
<u>Document ID: 15117</u>	1
<u>Introduction</u>	1
<u>Service Level Management Overview</u>	1
<u>Critical Success Factors</u>	1
<u>Performance Indicators</u>	2
<u>Service-level Management Process Flow</u>	2
<u>Implementing Service-level Management</u>	3
<u>Defining Network Service Levels</u>	3
<u>Creating and Maintaining SLAs</u>	23
<u>Service Level Management Performance Indicators</u>	30
<u>Documented Service Level Agreement or Service Level Definition</u>	31
<u>Performance Indicator Metrics</u>	32
<u>Service Level Management Review</u>	34
<u>Service Level Management Summary</u>	34
<u>Related Information</u>	34

Service Level Management: Best Practices White Paper

Document ID: 15117

Introduction

Service Level Management Overview

- Critical Success Factors

- Performance Indicators

- Service-level Management Process Flow

Implementing Service-level Management

- Defining Network Service Levels

- Creating and Maintaining SLAs

Service Level Management Performance Indicators

- Documented Service Level Agreement or Service Level Definition

- Performance Indicator Metrics

- Service Level Management Review

Service Level Management Summary

Related Information

Introduction

This document describes service-level management and service-level agreements (SLAs) for high-availability networks. It includes critical success factors for service-level management and performance indicators to help evaluate success. The document also provides significant detail for SLAs that follow best practice guidelines identified by the high availability service team.

Service Level Management Overview

Network organizations have historically met expanding network requirements by building solid network infrastructures and working reactively to handle individual service issues. When an outage occurred, the organization would build new processes, management capabilities, or infrastructure that to prevent a particular outage from occurring again. However, due to a higher change rate and increasing availability requirements, we now need an improved model to proactively prevent unplanned downtime and quickly repair the network. Many service-provider and enterprise organizations have attempted to better define the level of service required to achieve business goals.

Critical Success Factors

Critical success factors for SLAs are used to define key elements for successfully building obtainable service levels and for maintaining SLAs. To qualify as a critical success factor, a process or process step must improve the quality of the SLA and benefit network availability in general. The critical success factor should also be measurable so the organization can determine how successful it has been relative to the defined procedure.

See [Implementing Service-level Management](#) for more details.

Performance Indicators

Performance indicators provide the mechanism by which an organization measures critical success factors. You typically review these on a monthly basis to ensure that service-level definitions or SLAs are working well. The network operations group and the necessary tools groups can perform the following metrics.

Note: For organizations without SLAs, we recommend you perform service-level definitions and service-level reviews in addition to metrics.

Performance indicators include:

- Documented service-level definition or SLA that includes availability, performance, reactive service response time, problem resolution goals, and problem escalation.
- Monthly networking service-level review meeting to review service-level compliance and implement improvements.
- Performance indicator metrics, including availability, performance, service response time by priority, time to resolve by priority, and other measurable SLA parameters.

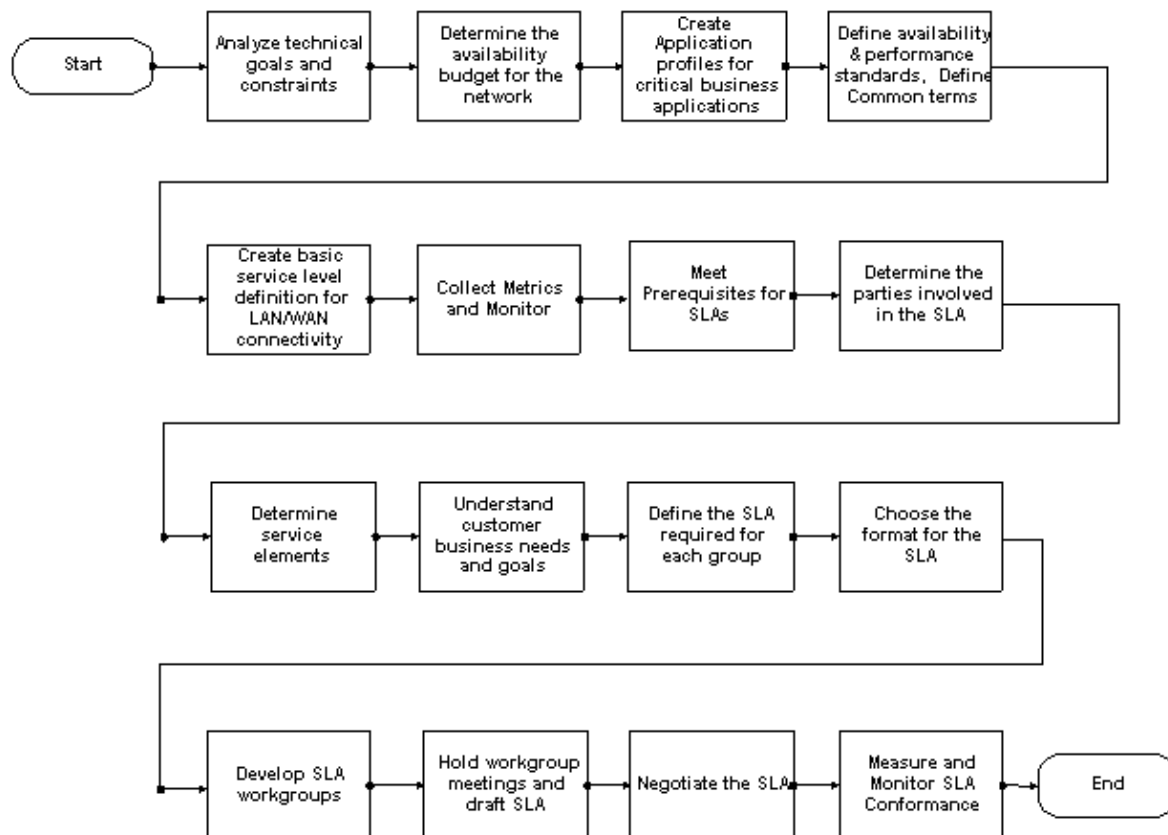
See *Implementing Service-level Management* for more information.

Service-level Management Process Flow

The high-level process flow for service-level management contains two major groups:

1. Defining network service levels
2. Creating and maintaining SLAs

Click on the objects in the following diagram to view the details for that step.



Implementing Service-level Management

Implementing service-level management consists of sixteen steps divided into the following two main categories:

- Defining network service levels
- Creating and maintaining SLAs

Defining Network Service Levels

Network managers need to define the major rules by which the network is supported, managed, and measured. Service levels provide goals for all network personnel and can be used as a metric in the quality of the overall service. You can also use service-level definitions as a tool for budgeting network resources and as evidence for the need to fund higher QoS. They also provide a way to evaluate vendor and carrier performance.

Without a service-level definition and measurement, the organization does not have clear goals. Service satisfaction may be governed by users with little differentiation between applications, server/client operations, or network support. Budgeting can be more difficult because the end result is not clear to the organization, and finally, the network organization tends to be more reactive, not proactive, in improving the network and support model.

We recommend the following steps for building and supporting a service-level model:

1. Analyze technical goals and constraints.
2. Determine the availability budget.
3. Create application profiles detailing network characteristics of critical applications.
4. Define availability and performance standards and define common terms.

5. Create a service–level definition that includes availability, performance, service response time, mean time to resolve problems, fault detection, upgrade thresholds, and escalation path.
6. Collect metrics and monitor the service–level definition.

Step 1: Analyze Technical Goals and Constraints

The best way to start analyzing technical goals and constraints is to brainstorm or research technical goals and requirements. Sometimes it helps to invite other IT technical counterparts into this discussion because these individuals have specific goals related to their services. Technical goals include availability levels, throughput, jitter, delay, response time, scalability requirements, new feature introductions, new application introductions, security, manageability, and even cost. The organization should then investigate constraints to achieving those goals given the available resources. You can create worksheets for each goal with an explanation of constraints. Initially, it may seem as if most of the goals are not achievable. Then start prioritizing the goals or lowering expectations that can still meet business requirements.

For example, you might have an availability level of 99.999 percent, or 5 minutes of downtime per year. There are numerous constraints to achieving this goal, such as single points of failure in hardware, mean time to repair (MTTR) broken hardware in remote locations, carrier reliability, proactive fault–detection capabilities, high change rates, and current network capacity limitations. As a result, you may adjust the goal to a more achievable level. The availability model in the next section can help you set realistic goals.

You may also think about providing higher availability in certain areas of the network that have fewer constraints. When the networking organization publishes service standards for availability, business groups within the organization may find the level unacceptable. This is then a natural point to begin SLA discussions or funding/budgeting models that can achieve the business requirements.

Work to identify all constraints or risks involved in achieving the technical goal. Prioritize constraints in terms of the greatest risk or impact to the desired goal. This helps the organization prioritize network improvement initiatives and determine how easily the constraint can be addressed. There are three kinds of constraints:

- Network technology, resiliency, and configuration
- Life–cycle practices, including planning, design, implementation, and operation
- Current traffic load or application behavior

Network technology, resiliency, and configuration constraints are any limitations or risks associated with the current technology, hardware, links, design, or configuration. Technology limitations cover any constraint posed by the technology itself. For example, no current technology allows for sub–second convergence times in redundant network environments, which may be critical for sustaining voice connections across the network. Another example may be the raw speed that data can traverse on terrestrial links, which is approximately 100 miles per millisecond.

Network hardware resiliency risk investigations should concentrate on hardware topology, hierarchy, modularity, redundancy, and MTBF along defined paths in the network. Network link constraints should focus on network links and carrier connectivity for enterprise organizations. Link constraints may include link redundancy and diversity, media limitations, wiring infrastructures, local–loop connectivity, and long–distance connectivity. Design constraints relate to the physical or logical design of the network and include everything from available space for equipment to scalability of the routing protocol implementation. All protocol and media designs should be considered in relation to configuration, availability, scalability, performance, and capacity. Network service constraints such as Dynamic Host Configuration Protocol (DHCP), Domain Name System (DNS), firewalls, protocol translators, and network address translators should also be considered.

Life-cycle practices define the processes and management of the network used to consistently deploy solutions, detect and repair problems, prevent capacity or performance problems, and configure the network for consistency and modularity. You need to consider this area because expertise and process are typically the largest contributors to non-availability. The network life cycle refers to the cycle of planning, design, implementation, and operations. Within each of these areas, you must understand network management functionality such as performance management, configuration management, fault management, and security. A network life-cycle assessment is available from Cisco NSA high-availability services (HAS) services showing current network availability constraints associated with network life-cycle practices.

Current traffic load or application constraints simply refer to the impact of current traffic and applications.

Unfortunately, many applications have significant constraints that require careful management. Jitter, delay, throughput, and bandwidth requirements for current applications typically have many constraints. The way the application was written may also create constraints. Application profiling helps you better understand these issues; the next section covers this feature. Investigating current availability, traffic, capacity, and performance overall also helps network managers to understand current service-level expectations and risks. This is typically accomplished with a process called network baselining, which helps to define network performance, availability, or capacity averages for a defined time period, normally about one month. This information is normally used for capacity planning and trending, but can also be used to understand service-level issues.

The following worksheet uses the above goal/constraint method for the example goal of preventing a security attack or denial-of-service (DoS) attack. You can also use this worksheet to help determine service coverage for minimizing security attacks.

Risk or Constraint	Type of Constraint	Potential Impact
Available DoS detection tools cannot detect all types of DoS attacks.	Technology/resiliency	High
Don't have the required staff and process to react to alerts.	Life-cycle practices	High
Current network access policies are not in place.	Life-cycle practices	Medium
Current lower-bandwidth Internet connection may be a factor if bandwidth congestion is used for attack.	Network capacity	Medium
Currently security configuration to help prevent attacks may not be thorough.	Technology/resiliency	Medium

Step 2: Determine the Availability Budget

An availability budget is the expected theoretical availability of the network between two defined points. Accurate theoretical information is useful in several ways:

- The organization can use this as a goal for internal availability and deviations can be quickly defined and remedied.

- The information can be used by network planners in determining the availability of the system to help ensure the design will meet business requirements.

Factors that contribute to non-availability or outage time include hardware failure, software failure, power and environmental issues, link or carrier failure, network design, human error, or lack of process. You should closely evaluate each of these parameters when evaluating the overall availability budget for the network.

If the organization currently measures availability, you may not need an availability budget. Use the availability measurement as a baseline to estimate the current service level used for a service-level definition. However, you may be interested in comparing the two to understand potential theoretical availability compared to the actual measured result.

Availability is the probability that a product or service will operate when needed. See the following definitions:

1. Availability

- ◆ $1 - (\text{total connection outage time}) / (\text{total in-service connection time})$
- ◆ $1 - [\text{Sigma}(\text{num connections affected in outage } i \times \text{duration of outage } i)] / (\text{num conns in service } \times \text{operating time})$

2. Unavailability

1 – Availability, or total outage connection time due to (hardware failure, software failure, environmental and power issues, link or carrier failure, network design, or user error and process failure)

3. Hardware Availability

The first area to investigate is potential hardware failure and the effect on unavailability. To determine this, the organization needs to understand the MTBF of all network components and the MTTR for hardware problems for all devices in a path between two points. If the network is modular and hierarchical, the hardware availability will be the same between almost any two points. MTBF information is available for all Cisco components and is available upon request to a local account manager. The Cisco NSA HAS program also uses a tool to help determine hardware availability along network paths, even when module redundancy, chassis redundancy, and path redundancy exist in the system. One major factor of hardware reliability is the MTTR. Organizations should evaluate how quickly they can repair broken hardware. If the organization has no sparing plan and relies on a standard Cisco SMARTnet" agreement, then the potential average replacement time is approximately 24 hours. In a typical LAN environment with core redundancy and no access redundancy, the approximate availability is 99.99 percent with a 4-hour MTTR.

4. Software Availability

The next area for investigation is software failures. For measurement purposes, Cisco defines software failures as device coldstarts due to software error. Cisco has made significant progress toward understanding software availability; however, newer releases take time to measure and are considered less available than general deployment software. General deployment software, such as IOS version 11.2(18), has been measured at over 99.9999 percent availability. This is calculated based on actual coldstarts on Cisco routers using six minutes as the repair time (time for router to reload). Organizations with a variety of versions are expected to have slightly lower availability because of added complexity, interoperability, and increased troubleshooting times. Organizations with the latest software versions are expected to have higher non-availability. The distribution for the non-availability is also fairly wide, meaning that customers could experience either significant non-availability or availability close to a general deployment release.

5. Environmental and Power Availability

You must also consider environmental and power issues in availability. Environmental issues relate to the breakdown of cooling systems needed to keep equipment at a specified operating temperature. Many Cisco devices will simply shut down when they are considerably out of specification rather than risking damage to all hardware. For the purpose of an availability budget, power will be used because it is the leading cause of non-availability in this area.

Although power failures are an important aspect of determining network availability, this discussion is limited because theoretical power analysis cannot be accurately done. What an organization must evaluate is an approximate measurement of power availability to its devices based on experience in its geographic area, power backup capabilities, and process implemented to ensure consistent quality power to all devices.

For a conservative evaluation, we can say that an organization with backup generators, uninterruptible-power-supply (UPS) systems, and quality power implementation processes may experience six 9s of availability, or 99.9999 percent, whereas organizations without these systems may experience availability at 99.99 percent, or approximately 36 minutes of downtime annually. Of course you can adjust these values to more realistic values based on the organization's perception or actual data.

6. Link or Carrier Failure

Link and carrier failures are major factors concerning availability in WAN environments. Keep in mind that WAN environments are simply other networks that are subject to the same availability issues as the organization's network, including hardware failure, software failure, user error, and power failure.

Many carrier networks have already performed an availability budget on their systems, but getting this information may be difficult. Keep in mind that carriers also frequently have availability guarantee levels that have little or no basis on an actual availability budget. These guarantee levels are sometimes simply marketing and sales methods used to promote the carrier. In some cases, these networks also publish availability statistics that appear extremely good. Keep in mind that these statistics may apply only to completely redundant core networks and don't factor in non-availability due to local-loop access, which is a major contributor to non-availability in WAN networks.

Creating an estimate of availability for WAN environments should be based on actual carrier information and the level of redundancy for WAN connectivity. If an organization has multiple building entrance facilities, redundant local-loop providers, Synchronous-Optical-Network (SONET) local access, and redundant long-distance carriers with geographic diversity, WAN availability will be considerably enhanced.

The phone service is a fairly accurate availability budget for non-redundant network connectivity in WAN environments. End-to-end connectivity for phones has an approximate availability budget of 99.94 percent using an availability budget methodology similar to the one described in this section. This methodology has been used successfully in data environments with only slight variation, and currently is being used as a target in the packet cable specification for service-provider cable networks. If we apply this value to a completely redundant system, we can assume that WAN availability will be close to 99.9999-percent available. Of course very few organizations have completely redundant, geographically dispersed WAN systems because of the expense and availability, so use proper judgement regarding this capability.

Link failures in a LAN environment are less likely. However, planners may want to assume a small

amount of downtime due to broken or loose connectors. For LAN networks, a conservative estimate is approximately 99.9999-percent availability, or about 30 seconds per year.

7. Network Design

Network design is another major contributor to availability. Non-scalable designs, design errors, and network convergence time all negatively affect availability.

Note: For the purposes of this document, non-scalable design or design errors are included in the following section.

Network design is then limited to a measurable value based on software and hardware failure in the network causing traffic re-routing. This value is typically called "system switchover time" and is a factor of the self-healing protocol capabilities within the system.

Calculate availability by simply using the same methods for system calculations. However, this is not valid unless the network switchover time meets network application requirements. If switchover time is acceptable, remove it from the calculation. If switchover time is not acceptable, then you must add it to the calculations. An example might be voice over IP (VoIP) in an environment where the estimated or actual switchover time is 30 seconds. In this example, users will simply hang up the phone and possibly try again. Users will certainly see this period of time as non-availability, yet it has not been estimated in the availability budget.

Calculate non-availability due to system switchover time by looking at the theoretical software and hardware availability along redundant paths, because switchover will occur in this area. You must know the number of devices that can fail and cause switchover in the redundant path, the MTBF of those devices, and the switchover time. A simple example would be a MTBF of 35,433 hours for each of two redundant identical devices and a switchover time of 30 seconds. Dividing 35,433 by 8766 (hours per year averaged to include leap years), we see that the device will fail once every four years. If we use 30 seconds as a switchover time, we can then assume that each device will experience, on average, 7.5 seconds per year of non-availability due to switchover. Since users may be traversing either path, the result is then doubled to 15 seconds per year. When this is calculated in terms of seconds per year, the amount of availability due to switchover can be calculated as 99.99999785-percent availability in this simple system. This may be higher in other environments because of the number of redundant devices in the network where switchover is a potential.

8. User Error and Process

User error and process availability issues are the major causes of non-availability in enterprise and carrier networks. Approximately 80 percent of non-availability occurs because of issues such as not detecting errors, change failures, and performance problems.

Organizations will simply not want to use four times all other theoretical non-availability in determining the availability budget, yet evidence consistently suggests that this is the case in many environments. The next section covers this aspect of non-availability more thoroughly.

Since you cannot theoretically calculate the amount of non-availability due to user error and process, we recommend you remove this removed from the availability budget and that organizations strive for perfection. The one caveat is that organizations need to understand the current risk to availability in their own processes and levels of expertise. Once you better understand these risks and inhibitors, network planners may wish to factor in some quantity of non-availability due to these issues. The Cisco NSA HAS program investigates these issues and can help organizations understand potential non-availability due to process, user error, or expertise issues.

9. Determining the Final Availability Budget

You can determine the overall availability budget by multiplying availability for each of the previously defined areas. This is typically done for homogenous environments where connectivity is similar between any two points, such as a hierarchical modular LAN environment or a hierarchical standard WAN environment.

In this example, the availability budget is done for a hierarchical modular LAN environment. The environment uses backup generators and UPS systems for all network components and properly manages power. The organization does not use VoIP and does not wish to factor in software switchover time. The estimates are:

- ◆ Hardware path availability between two end points = 99.99 percent availability
- ◆ Software availability using GD software reliability as reference = 99.9999 percent availability
- ◆ Environmental and power availability with backup systems = 99.999 percent availability
- ◆ Link failure in LAN environment = 99.9999 percent availability
- ◆ System switchover time not factored = 100 percent availability
- ◆ User error and process availability assumed perfect = 100 percent availability

The final availability budget that the organizations should strive for equals $0.9999 \times 0.999999 \times 0.999999 \times 0.999999 = 0.999896$, or 99.9896 percent availability. If we factor in potential non-availability due to user or process error and assume that non-availability is 4X availability due to technical factors, we could assume that the availability budget is 99.95 percent.

This example analysis indicates then that LAN availability would fall on average between 99.95 and 99.989 percent. These numbers can now be used as a service level goal for the networking organization. You can gain additional value by measuring availability in the system and determining what percentage of non-availability was due to each of the above six areas. This allows the organization to properly evaluate vendors, carriers, processes, and staff. The number can also be used to set expectations within the business. If the number is unacceptable, then budget additional resources to gain the desired levels.

It may be useful for network managers to understand the amount of downtime at any particular availability level. The amount of downtime in minutes for a one-year period, given any availability level, is:

Minutes of downtime in one year = $525600 - (\text{Availability level} \times 5256)$

If you use the availability level of 99.95 percent, this works out to be equal to $525600 - (99.95 \times 5256)$, or 222.8 minutes of downtime. For the above availability definition, this is equal to the average amount of downtime for all connections in service within the network.

Step 3: Create Application Profiles

Application profiles help the networking organization understand and define network service level requirements for individual applications. This helps to ensure that the network supports individual application requirements and network services overall. Application profiles can also serve as a documented baseline for network service support when application or server groups point to the network as the problem. Ultimately, application profiles help to align network service goals with application or business requirements by comparing application requirements such as performance and availability with realistic network service goals or current limitations. This is important not only for service level management, but also for overall top-down network design.

Create application profiles any time you introduce new applications to the network. You may need an agreement between the IT application group, server administration groups, and networking to help enforce

application profile creation for new and existing services. Complete application profiles for business applications and system applications. Business applications may include e-mail, file transfer, Web browsing, medical imaging, or manufacturing. System applications may include software distribution, user authentication, network backup, and network management.

A network analyst and an application or server support application should create the application profile. New applications may require the use of a protocol analyzer and WAN emulator with delay emulation to properly characterize application requirements. This helps identify the necessary bandwidth, maximum delay for application usability, and jitter requirements. This can be done in a lab environment as long as you have the required servers. In other cases, such as with VoIP, network requirements including jitter, delay, and bandwidth are well published and lab testing will not be needed. An application profile should include the following items:

- Application name
- Type of application
- New application?
- Business importance
- Availability requirements
- Protocols and ports used
- Estimated user bandwidth (kbps)
- Number and location of users
- File transfer requirements (including time, volume, and endpoints)
- Network outage impact
- Delay, jitter, and availability requirements

The goal of the application profile is to understand business requirements for the application, business criticality, and network requirements such as bandwidth, delay, and jitter. In addition, the networking organization should understand the impact of network downtime. In some cases, you will need application or server re-starts that significantly add to overall application downtime. When you complete the application profile, you can compare overall network capabilities and help to align network service levels with business and application requirements.

Step 4: Define Availability and Performance Standards

Availability and performance standards set the service expectations for the organization. These may be defined for different areas of the network or specific applications. Performance may also be defined in terms of round-trip delay, jitter, maximum throughput, bandwidth commitments, and overall scalability. In addition to setting the service expectations, the organization should also take care to define each of the service standards so that user and IT groups working with networking fully understand the service standard and how it relates to their application or server administration requirements. User and IT groups should also understand how the service standard might be measured.

Results from previous service level definition steps will help to create the standard. At this point, the networking organization should have a clear understanding of the current risks and constraints in the network, an understanding of application behavior, and a theoretical availability analysis or availability baseline.

1. Define the geographic or application areas where service standards will be applied.

This may include areas such as the campus LAN, domestic WAN, extranet, or partner connectivity. In some cases, the organization may have different service level goals within one area. This is not uncommon for enterprise or service provider organizations. In these cases, it would not be uncommon to create different service level standards based on individual service requirements. These may be

classified as gold, silver, and bronze service standards within one geographic or service area.

2. Define the service standard parameters.

Availability and round-trip delay are the most common network service standards. Maximum throughput, minimum bandwidth commitment, jitter, acceptable error rates, and scalability capabilities may also be included as needed. Be careful when reviewing the service parameter for measurement methods. Whether or not the parameter moves on to a SLA, the organization should think about how the service parameter might be measured or justified when problems or service disagreements occur.

After you define the service areas and service parameters, use the information from previous steps to build a matrix of service standards. The organization will also need to define areas that may be confusing to users and IT groups. For instance, the maximum response time will be very different for a round-trip ping than for hitting the Enter key at a remote location for a specific application. The following table shows the performance targets within the United States.

Network Area	Availability Target	Measurement Method	Average Network Response Time Target	Max Response Time Accepted	Response Time Measurement Method
LAN	99.99%	Impacted user minutes	Under 5 ms	10 ms	Round-trip ping response
WAN	99.99%	Impacted user minutes	Under 100 ms (round-trip ping)	150 ms	Round-trip ping response
Critical WAN and Extranet	99.99%	Impacted user minutes	Under 100 ms (round-trip ping)	150 ms	Round-trip ping response

Step 5: Define Network Service

This is the last step toward basic service level management; it defines the reactive and proactive processes and network management capabilities that you implement to achieve service level goals. The final document is typically called an operations support plan. Most application support plans include only reactive support requirements. In high-availability environments, the organization must also consider proactive management processes that will be used to isolate and resolve network issues before user service calls are initiated. Overall, the final document should:

- Describe the reactive and proactive process used to achieve the service level goal
- How the service process will be managed
- How the service goal and service process will be measured.

This section contains examples for reactive service definitions and proactive service definitions to consider for many service-provider and enterprise organizations. The goal in building the service level definitions is to create a service that will meet the availability and performance goals. To accomplish this, the organization must build the service with the current technical constraints, availability budget, and application profiles in mind. Specifically, the organization should define and build a service that consistently and quickly identifies

and resolves problems within times allocated by the availability model. The organization must also define a service that can quickly identify and resolve potential service issues that will impact availability and performance if ignored.

You will not achieve the desired service level overnight. Shortcomings such as low expertise, current process limitations, or inadequate staffing levels may prevent the organization from achieving the desired standards or goals, even after the previous service analysis steps. There is no precise method to exactly match the required service level to the desired goals. To accommodate for this, the organization should measure the service standards and measure the service parameters used to support the service standards. When the organization is not meeting service goals, it should then look to service metrics to help understand the issue. In many cases, budgeting increases can be made to improve support services and make improvements necessary to achieve the desired service goals. Over time the organization may make several adjustments, either to the service goal or to the service definition, in order to align network services and business requirements.

For example, an organization might achieve 99 percent availability when the goal was much higher at 99.9 percent availability. When looking at service and support metrics, representatives of the organization found that hardware replacement was taking approximately 24 hours, much longer than the original estimate because the organization had budgeted only four. In addition, the organization found that proactive management capabilities were being ignored and down redundant network devices were not being repaired. They also found that they didn't have the personnel to make improvements. As a result, after considering lowering the current service goals, the organization budgeted for additional resources needed to achieve the desired service level.

Service definitions should include both reactive support definitions and proactive definitions. Reactive definitions define how the organization will react to problems after they have been identified from either user complaint or network management capabilities. Proactive definitions describe how the organization will identify and resolve potential network problems, including repair of broken "standby" network components, error detection, and capacity thresholds and upgrades. The following sections provide examples of both reactive and proactive service level definitions.

Reactive Service Level Definitions

The following service level areas are typically measured using help–desk database statistics and periodic auditing. This table shows example of problem severity for an organization. Notice that the chart does not include how to handle requests for new service, which may be handled by a SLA or additional application profiling and performance what–if analysis. Typically severity 5 may be a request for new service if handled via the same support process.

Severity 1	Severity 2	Severity 3	Severity 4
Severe business impact	High business impact through loss or degradation, possible workaround in place	Some specific network functionality is lost or degraded, such as loss of redundancy	A functional query or fault that has no business impact for the organization
LAN user or server segment down	Campus LAN down; 5–99 users affected	Campus LAN performance impacted	
Critical WAN site down			

Domestic WAN site down	LAN redundancy lost	
International WAN site down		
Critical performance impact		

When problem severity has been defined, define or investigate the support process to create service response definitions. In general, service response definitions require a tiered support structure coupled with a help desk software support system to track problems via trouble tickets. Metrics should also be available on response time and resolution time for each priority, number of calls by priority, and response/resolution quality. To define the support process, it helps to define the goals of each support tier in the organization and their roles and responsibilities. This helps the organization understand resource requirements and levels of expertise for each support level. The following table provides an example of a tiered support organization with problem resolution guidelines.

Support Tier	Responsibility	Goals
Tier 1 Support	Full-time help desk support Answer support calls, place trouble tickets, work on problem up to 15 minutes, document ticket and escalate to appropriate tier 2 support	Resolution of 40% of incoming calls
Tier 2 Support	Queue monitoring, network management, station monitoring Place trouble tickets for software identified problems Implement Take calls from tier 1, vendor, and tier 3 escalation Assume ownership of call until resolution	Resolution of 100% of calls at tier 2 level
Tier 3 Support	Must provide immediate support to tier 2 for all priority 1 problems Agree to help with all problems unsolved by tier 2 within SLA resolution period	No direct problem ownership

The next step is to create the matrix for the service response and service resolution service definition. This sets goals for how quickly problems are resolved, including hardware replacement. It is important to set goals in this area because service response time and recovery time directly impact network availability. Problem

resolution times should also be aligned with the availability budget. If large numbers of high severity problems are not accounted for in the availability budget, the organization can then work to understand the source of these problems and a potential remedy. See the following table:

Problem Severity	Help Desk	Tier 2	Onsite	Hardware	Problem
1	Response Immediate escalation to tier 2, network operations manager	Response 5 minutes	Tier 2 2 hours	Replacement 2 hours	Resolution 4 hours
2	Immediate escalation to tier 2, network operations manager	5 minutes	4 hours	4 hours	8 hours
3	15 minutes	2 hours	12 hours	24 hours	36 hours
4	15 minutes	4 hours	3 days	3 days	6 days

In addition to service response and service resolution, build a matrix for escalation. The escalation matrix helps ensure that available resources are focused on problems that severely affect service. In general, when analysts are focused on fixing problems, they rarely focus on bringing additional resources in on the problem. Defining when additional resources should be notified helps to promote problem awareness in management and can generally help lead to future proactive or preventative measures. See the following table:

Elapsed Time	Severity 1	Severity 2	Severity 3	Severity 4
5 minutes	Network operations manager, tier 3 support, director of networking			
1 hour	Update to network operations manager, tier 3 support, director of networking	Update to network operations manager, tier 3 support,		
2 hours	Escalate to VP, update to director, operations manager	director of networking		

4 hours	Root cause analysis to VP, director, operations manager, tier 3 support, unresolved requires CEO notification	Escalate to VP, update to director, operations manager		
24 hours			Network operations manager	
5 days				Network operations manager

So far, the service level definitions have focused on how the operations support organization reacts to problems after they are identified. Operations organizations have created operational support plans with information similar to the above for years. However, what is missing in these cases is how the organization will identify problems and which problems they will identify. More sophisticated network organizations have attempted to resolve this issue by simply creating goals for the percentage of problems that are proactively identified, as opposed to problems reactively identified by user problem report or complaint.

The next table shows how an organization may wish to measure proactive support capabilities and proactive support overall.

Network Area	Proactive Problem	Reactive Problem
LAN	Identification Ratio 80 %	Identification Ratio 20 %
WAN	80 %	20 %

This is a good start at defining more proactive support definitions because it is simple and fairly easy to measure, especially if proactive tools automatically generate trouble tickets. This also helps focus network management tools/information on resolving problems proactively rather than helping with the root cause. However, the main issue with this method is that it does not define proactive support requirements. This generally creates gaps in proactive support management capabilities and results in additional availability risk.

Proactive Service Level Definitions

A more comprehensive methodology for creating service level definitions includes more detail on how the network is monitored and how the operations organization reacts to defined network management station (NMS) thresholds on a 7 x 24 basis. This may seem like an impossible task given the sheer number of Management Information Base (MIB) variables and the amount of network management information available that is pertinent to network health. It could also be extremely expensive and resource intensive. Unfortunately, these objections prevent many from implementing a proactive service definition that, by nature, should be simple, fairly easy to follow, and applicable only to the greatest availability or performance risks in the network. If an organization then sees value in basic proactive service definitions, more variables can be added over time without significant impact, as long as you implement a phased approach.

Include the first area of proactive service definitions in all operations support plans. The service definition simply states how the operations group will proactively identify and respond to network or link down conditions in different areas of the network. Without this definition (or management support), the organization can expect variable support, unrealistic user expectations, and ultimately lower network availability.

The following table shows how an organization might create a service definition for link/device–down conditions. The example shows an enterprise organization that may have different notification and response requirements based on the time of day and area of the network.

Network Device or Link Down	Detection	5 x 8	7 x 24	5 x 8	7 x 24
Core LAN	Method SNMP device and link polling, traps	Notification NOC creates trouble ticket, page LAN–duty pager	Notification Auto page LAN duty pager, LAN duty person creates trouble ticket for core LAN queue	Resolution LAN analyst assigned within 15 minutes by NOC, repair as per service response definition	Resolution Priorities 1 and 2 immediate investigation and resolution Priorities 3 and 4 queue for morning resolution
Domestic WAN	SNMP device and link polling, traps	NOC creates trouble ticket, page WAN duty pager	Auto page WAN duty pager, WAN duty person creates trouble ticket for WAN queue	WAN analyst assigned within 15 minutes by NOC, repair as per service response definition	Priorities 1 and 2 immediate investigation and resolution Priorities 3 and 4 queue for morning resolution
Extranet	SNMP device and link polling, traps	NOC creates trouble ticket, page partner duty pager	Auto page partner duty pager, partner duty person creates trouble ticket for partner queue	Partner analyst assigned within 15 minutes by NOC, repair as per service response definition	Priorities 1 and 2 immediate investigation and resolution; Priorities 3 and 4 queue for morning resolution

The remaining proactive service level definitions can be divided into two categories: network errors and capacity/performance issues. Only a small percentage of network organizations have service level definitions in these areas. As a result, these issues are ignored or handled sporadically. This may be fine in some network

environments, but high availability environments will generally require consistent proactive service management.

Networking organizations tend to struggle with proactive service definitions for several reasons. This is primarily because they have not performed a requirements analysis for proactive service definitions based on availability risks, the availability budget, and application issues. This leads to unclear requirements for proactive service definitions and unclear benefits, especially because additional resources may be needed.

The second reason involves balancing the amount of proactive management that can be done with existing or newly-defined resources. Only generate those alerts that have serious potential impact to availability or performance. You must also consider event correlation management or processes to ensure that multiple proactive trouble tickets are not generated for the same problem. The last reason organizations may struggle is that creating a new set of proactive alerts can often generate an initial flood of messages that have previously gone undetected. The operations group must be prepared for this initial flood of issues and additional short-term resources to fix or resolve these previously undetected conditions.

The first category of proactive service level definitions is network errors. Network errors can be further subdivided into system errors that include software errors or hardware errors, protocol errors, media control errors, accuracy errors, and environmental warnings. Developing a service level definition starts with a general understanding of how these problem conditions will be detected, who will look at them, and what will happen when they occur. Add specific messages or issues to the service level definition if the need arises. You may also need additional work in the following areas to ensure success:

- Tier 1, tier 2, and tier 3 support responsibilities
- Balancing the priority of the network management information with the amount of proactive work that the operations group can effectively handle
- Training requirements to ensure support staff can effectively deal with the defined alerts
- Event correlation methodologies to ensure that multiple trouble tickets are not generated for the same root-cause problem
- Documentation on specific messages or alerts that helps with event identification at tier 1 support level

The following table shows an example service level definition for network errors that provide a clear understanding of who is responsible for proactive network error alerts, how the problem will be identified, and what will happen when the problem occurs. The organization may still need additional efforts as defined above to ensure success

s.

Error Category	Detection Method	Threshold	Action Taken
Software Errors (crashes forced by software)	Daily review of syslog messages using syslog viewer Done by tier 2 support	Any occurrence for priority 0, 1, and 2 Over 100 occurrences of level 3 or above	Review problem, create trouble ticket, and dispatch if new occurrence or if problem requires attention

Hardware Errors (crashes forced by hardware)	Daily review of syslog messages using syslog viewer Done by tier 2 support	Any occurrence for priority 0, 1, and 2 Over 100 occurrences of level 3 or above	Review problem, create trouble ticket, and dispatch if new occurrence or if problem requires attention
Protocol Errors (IP routing protocols only)	Daily review of syslog messages using syslog viewer Done by tier 2 support	Ten messages per day of priorities 0, 1, and 2 Over 100 occurrences of level 3 or above	Review problem, create trouble ticket, and dispatch if new occurrence or if problem requires attention
Media Control Errors (FDDI, POS, and Fast Ethernet only)	Daily review of syslog messages using syslog viewer Done by tier 2 support	Ten messages per day of priorities 0, 1, and 2 Over 100 occurrences of level 3 or above	Review problem, create trouble ticket, and dispatch if new occurrence or if problem requires attention
Environmental Messages (power and temp)	Daily review of syslog messages using syslog viewer Done by tier 2 support	Any message	Create trouble ticket and dispatch for new problems
Accuracy Errors (link input errors)	SNMP polling at 5-minute intervals Threshold events received by NOC	Input or output errors One error in any 5-minute interval on any link	Create trouble ticket for new problems and dispatch to tier 2 support

The other category of proactive service level definitions applies to performance and capacity. True performance and capacity management includes exception management, baselining and trending, and what-if analysis. The service level definition simply defines performance and capacity exception thresholds and average thresholds that will initiate investigation or upgrade. These thresholds may then apply to all three

performance and capacity management processes in some way.

Capacity and performance service level definitions can be broken down into several categories: network links, network devices, end-to-end performance, and application performance. Developing service level definitions in these areas requires in-depth technical knowledge regarding specific aspects of device capacity, media capacity, QoS characteristics, and application requirements. For this reason, we recommend that network architects develop performance and capacity-related service level definitions with vendor input.

Like network errors, developing a service level definition for capacity and performance starts with a general understanding of how these problem conditions will be detected, who will look at them, and what will happen when they occur. You can add specific event definitions to the service level definition if the need arises. You may also need additional work in the following areas to ensure success:

- A clear understanding of application performance requirements
- In-depth technical investigation on threshold values that make sense for the organization based on business requirements and overall costs
- Budgetary cycle and out-of-cycle upgrade requirements
- Tier 1, tier 2, and tier 3 support responsibilities
- Priority and criticality of the network management information balanced with the amount of proactive work that the operations group can effectively handle
- Training requirements to ensure that support staff understand the messages or alerts and can effectively deal with the defined condition
- Event correlation methodologies or processes to ensure that multiple trouble tickets are not generated for the same root-cause problem
- Documentation on specific messages or alerts that helps with event identification at the tier 1 support level

The following table shows an example service level definition for link utilization that provides a clear understanding of who is responsible for proactive network error alerts, how the problem will be identified, and what will happen when the problem occurs. The organization may still need additional efforts as defined above to ensure success.

Network Area/Media	Detection Method	Threshold	Action Taken
Campus LAN Backbone and Distribution Links	SNMP polling at 5-minute intervals	50% utilization in 5-minute intervals	E-mail notification to performance e-mail alias
	RMON exception traps on core and distribution links	90% utilization via exception	Group to evaluate QoS requirement or plan upgrade for
Domestic WAN Links	SNMP polling at 5-minute intervals	75% utilization in 5-minute intervals	recurring issues E-mail notification to performance e-mail alias Group to evaluate QoS

			requirement or plan upgrade for recurring issues
Extranet WAN Links	SNMP polling at 5-minute intervals	60% utilization in 5-minute intervals	E-mail notification to performance e-mail alias Group to evaluate QoS requirement or plan upgrade for recurring issues

The following table defines service level definitions for device capacity and performance thresholds. Ensure you create thresholds that are meaningful and useful in preventing network problems or availability issues. This is a very important area because un-checked device control plane resource issues can have serious network impact.

Cisco 7500	CPU, memory, buffers	SNMP polling at 5-minute intervals RMON notification for CPU	CPU at 75% during 5-minute intervals, 99% via RMON notification Memory at 50% during 5-minute intervals Buffers at 99% utilization	E-mail notification to performance and capacity e-mail alias group to resolve issues or plan upgrade RMON CPU at 99%, place trouble ticket and page tier 2 support pager
Cisco 2600	CPU, memory	SNMP polling at 5-minute intervals	CPU at 75% during 5-minute intervals Memory at 50% during 5-minute intervals	E-mail notification to performance and capacity e-mail alias group to resolve issues or
Catalyst 5000	Backplane utilization,	SNMP polling at	Backplane at 50%	plan upgrade E-mail notification

	memory	5-minute intervals	utilization Memory at 75% utilization	to performance and capacity e-mail alias group to resolve issues or plan upgrade
LightStream® 1010 ATM switch	CPU, memory	SNMP polling at 5-minute intervals	CPU at 65% utilization Memory at 50% utilization	E-mail notification to performance and capacity e-mail alias group to resolve issues or plan upgrade

The next table defines service level definitions for end-to-end performance and capacity. These thresholds are generally based on application requirements but can also be used to indicate some type of network performance or capacity problem. Most organizations with service level definitions for performance create only a handful of performance definitions because measuring performance from every point in the network to every other point requires significant resources and creates a high amount of network overhead. These end-to-end performance issues may also be caught in link or device capacity thresholds. We recommend general definitions by geographic area. Some critical sites or links may be added if necessary.

Network Area/Media	Measurement	Threshold	Action Taken
Campus LAN	Method None		
	No problem expected Difficult to measure entire LAN infrastructure	10-millisecond round-trip response time or less at all times	E-mail notification to performance and capacity e-mail alias group to resolve issue
Domestic WAN Links	Current measurement from SF to NY and SF to Chicago only using Internet Performance Monitor (IPM) ICMP echo	75-millisecond round-trip response time averaged over 5-minute period	or plan upgrade E-mail notification to performance e-mail alias group to evaluate QoS requirement or plan upgrade for recurring issues

San Francisco to Tokyo	Current measurement from San Francisco to Brussels using IPM and ICMP echo	250–millisecond round–trip response time averaged over 5–minute period	E–mail notification to performance e–mail alias group to evaluate QoS requirement or plan upgrade for recurring issues
San Francisco to Brussels	Current measurement from San Francisco to Brussels using IPM and ICMP echo	175–millisecond round–trip response time averaged over 5–minute period	E–mail notification to performance e–mail alias group to evaluate QoS requirement or plan upgrade for recurring issues

The final area for service level definitions is for application performance. Application performance service level definitions are normally created by the application or server administration group because performance and capacity of the servers themselves is probably the largest factor in application performance. Networking organizations can realize tremendous benefit by creating service level definitions for network application performance because:

- service level definitions and measurement can help eliminate conflicts between groups.
- service level definitions for individual applications are important if QoS is configured for key applications and other traffic is considered optional.

If you choose to create and measure application performance, it is probably best if you do not measure performance to the server itself. This then helps distinguish between network problems and application or server problems. Use probes or the system availability agent software running on Cisco routers and the Cisco IPM controlling the packet type and measurement frequency.

The following table shows a simple service level definition for application performance.

Application	Measurement Method	Threshold	Action Taken
Enterprise Resource Planning (ERP) Application TCP Port 1529 Brussels to	Brussels to San Francisco using IPM measuring port 1529 round–trip performance	175–millisecond round–trip response time averaged over 5–minute period	E–mail notification to performance e–mail alias group to evaluate problem or plan upgrade for recurring

SF	Brussels gateway to SFO gateway 2		issues
ERP Application TCP Port 1529 Tokyo to SF	Brussels to San Francisco using IPM measuring port 1529 round-trip performance Brussels gateway to SFO gateway 2	200-millisecond round-trip response time averaged over 5-minute period	E-mail notification to performance e-mail alias group to evaluate problem or plan upgrade for recurring issues
Customer Support Application TCP port 1702 Sydney to SF	Sydney to San Francisco using IPM measuring port 1702 round-trip performance Sydney gateway to SFO gateway	250-millisecond round-trip response time averaged over 5-minute period	E-mail notification to performance e-mail alias group to evaluate problem or plan upgrade for recurring issues

1

Step 6: Collect Metrics and Monitor

service level definitions by themselves are worthless unless the organization collects metrics and monitors success. In creating a critical service level definition, define how the service level will be measured and reported. Measuring the service level determines whether the organization is meeting objectives and also identifies the root cause of availability or performance issues. Also consider the goal when choosing a method to measure the service level definition. See *Creating and Maintaining SLAs* for more information.

Monitoring service levels entails conducting a periodic review meeting, normally every month, to discuss periodic service. Discuss all metrics and whether they conform to the objectives. If they do not conform, determine the root cause of the problem and implement improvements. You should also cover current initiatives and progress in improving individual situations.

Creating and Maintaining SLAs

service level definitions are an excellent building block in that they help create a consistent QoS throughout the organization and help improve availability. The next step is SLAs, which are an improvement because they align business objectives and cost requirements directly to service quality. The well-constructed SLA then serves as a model for efficiency, quality, and synergy between the user community and support group by maintaining clear processes and procedures for network issues or problems.

SLAs provide several benefits:

- SLAs establish two-way accountability for service, meaning that users and application groups are also accountable for the network service. If they don't help create a SLA for a specific service and communicate business impact with the network group, then they may actually be accountable for the problem.
- SLAs help determine standard tools and resources needed to meet business requirements. Deciding how many people and which tools to use without SLAs is often a budgetary guess. The service may be over-engineered, which leads to over-spending, or under-engineered, which leads to unmet business objectives. Tuning SLAs helps achieve that balanced optimal level.
- The documented SLA creates a clearer vehicle for setting service level expectations.

We recommend the following steps for building SLAs after service level definitions have been created: We recommend the following steps for building SLAs after service level definitions have been created:

7. Meet the prerequisites for SLAs.
8. Determine the parties involved in the SLA.
9. Determine service elements.
10. Understand customer business needs and goals
11. Define the SLA required for each group.
12. Choose the format of the SLA
13. Develop SLA workgroups
14. Hold workgroup meetings and draft the SLA.
15. Negotiate the SLA.
16. Measure and monitor SLA conformance.

Step 7: Meet Prerequisites for SLAs

Experts in IT SLA development identified three prerequisites to a successful SLA. Unfortunately, organizations that do not meet these objectives can expect problems with the SLA process and should consider the potential problems involved with the SLA process. Failing to implement SLAs is not detrimental if the networking organization can build service level definitions that meet general business requirements. The following are prerequisites for the SLA process:

- Your business must have a service-oriented culture.

The organization must place the needs of the customers first. You need a top-down priority commitment to service, resulting in a complete understanding of customer needs and perceptions. Conduct customer satisfaction surveys and customer-driven service initiatives.

Another service indicator may be that the organization states service or support satisfaction as a corporate goal. This is not uncommon because IT organizations are now critically linked to overall organization success.

The service culture is important because the SLA process is fundamentally about making

improvements based on customer needs and business requirements. If organizations have not done this in the past, they will find the SLA process difficult.

- Customer/business initiatives must drive all IT activities.

The company vision or mission statements must be aligned with customer and business initiatives, which then drive all IT activities, including SLAs. Too often a network is put in place to meet a particular goal, yet the networking group loses sight of that goal and subsequent business requirements. In these cases, a set budget is allocated to the network, which may overreact to current needs or grossly underestimate the requirement, resulting in failure.

When customer/business initiatives are aligned with IT activities, the networking organization can more easily be in tune with new application rollouts, new services, or other business requirements. The relationship and common overall focus on meeting corporate goals are present and all groups execute as a team.

- You must commit to the SLA process and contract.

First there must be commitment to learn the SLA process to develop effective agreements. Second, you must honor the service requirements of the contract. Don't expect to create powerful SLAs without significant input and commitment from all individuals involved. This commitment must also come from management and all individuals associated with the SLA process.

Step 8: Determine the Parties Involved in the SLA

Enterprise-level network SLAs depend heavily on network elements, server administration elements, help-desk support, application elements, and business or user requirements. Normally management from each area will be involved in the SLA process. This scenario works well when the organization is building basic reactive support SLAs. Enterprise organizations with higher-availability requirements may need technical assistance during the SLA process to help with such issues as availability budgeting, performance limitations, application profiling, or proactive management capabilities. For more proactive management SLA aspects, we recommend a technical team of network architects and application architects. Technical assistance can much more closely approximate the availability and performance capabilities of the network and what would be needed to reach specific objectives.

Service-provider SLAs do not normally include user input because they are created for the sole purpose of gaining a competitive edge on other service providers. In some cases, upper management will create these SLAs at very high-availability or high-performance levels to promote their service and to provide internal goals for internal employees. Other service providers will concentrate on the technical aspects of improving availability by creating strong service level definitions that are measured and managed internally. In other cases, both efforts occur simultaneously but not necessarily together or with the same goals.

Choosing the parties involved in the SLA should then be based on the goals of the SLA. Some possible goals are:

- Meeting reactive support business objectives
- Providing the highest level of availability by defining proactive SLAs
- Promoting or selling a service

Step 9: Determine Service Elements

Primary service/support SLAs will normally have many components, including the level of support, how it will be measured, the escalation path for SLA reconciliation, and overall budget concerns. Service elements for high-availability environments should include proactive service definitions as well as reactive goals.

Additional details include the following:

- Onsite support business hours and procedures for off-hours support
- Priority definitions, including problem type, maximum time to begin work on the problem, maximum time to resolve the problem, and escalation procedures
- Products or services to be supported, ranked in order of business criticality
- Support for expertise expectations, performance-level expectations, status reporting, and user responsibilities for problem resolution
- Geographic or business unit support-level issues and requirements
- Problem management methodology and procedures (call-tracking system)
- Help desk goals
- Network error detection and service response
- Network availability measurement and reporting
- Network capacity and performance measurement and reporting
- Conflict resolution procedures
- Funding the implemented SLA

Networked application or service SLAs may have additional needs based on user group requirements and business criticality. The network organization must listen closely to these business requirements and develop specialized solutions that fit into the overall support structure. Fitting into the overall support culture is critical because it is important not to create a premier service intended only for some individuals or groups. In many cases, these additional requirements can be placed into "solution" categories. An example might be a platinum, gold, and silver solution based on business need. See the following examples of SLA requirements for specific business needs.

Note: The support structure, escalation path, help-desk procedures, measurement, and priority definitions should largely remain the same to maintain and improve a consistent service culture.

- Bandwidth requirements and capabilities for burst
- Performance requirements
- QoS requirements and definitions
- Availability requirements and redundancy to build solution matrix
- Monitoring and reporting requirements, methodology, and procedures
- Upgrade criteria for application/service elements
- Funding out-of-budget requirements or cross-charging methodology

For instance, you can create solution categories for WAN site connectivity. The platinum solution would be provided with twin T1 services to the site. A different carrier would provide each T1 line. The site would have two routers configured so that if any T1 or router failed the site would not experience an outage. The gold service would have two routers, but backup Frame Relay would be used. This solution may have limited bandwidth for the duration of the outage. The silver solution would have only one router and one carrier service. Any of these solutions would be considered for different priority levels for problem tickets. Some organizations may require a platinum or gold solution if a priority 1 or 2 ticket is required for an outage. Customer organizations can then fund the level of service they require. The following table shows an example of an organization that offers three levels of service, depending on business need for extranet connectivity.

Solution	Platinum	Gold	Silver
<i>Devices</i>	Redundant routers for WAN connectivity	Redundant router for backup at core site	No device redundancy
<i>WAN</i>			

	Redundant T1 connectivity, multiple carriers	T1 connectivity with Frame Relay backup	No WAN redundancy
<i>Bandwidth Requirements and Burst</i>	Redundant T1 with load sharing for burst	Non-load sharing, Frame Relay backup for critical applications only; Frame Relay 64K CIR only	UP to T1
<i>Performance</i>	Consistent 100-ms round-trip response time or less	Response time 100 ms or less expected 99.9%	Response time 100 ms or less expected
<i>Availability Requirement</i>	99.99%	99.95%	99%
<i>Help desk Priority when Down</i>	Priority 1: business-critical service down	Priority 2: business-impacting service down	Priority 3: business connectivity down

Step 10: Understand Customer Business Needs and Goals

This step lends the SLA developer a great deal of credibility. By understanding the needs of the various business groups, the initial SLA document will be much closer to the business requirement and desired result. Try to understand the cost of downtime for the customer's service. Estimate in terms of lost productivity, revenue, and customer goodwill. Keep in mind that even simple connections with a few people can seriously impact revenue. In this case, be sure to help the customer understand the availability and performance risks that may occur so that the organization better understands the level of service it needs. If you miss this step, you may get many customers simply demanding 100-percent availability.

The SLA developer should also understand the business goals and growth of the organization in order to accommodate network upgrades, workload, and budgeting. It is also helpful to understand the applications that will be used. Hopefully the organization has application profiles on each application, but if not, consider doing a technical evaluation of the application to determine network-related issues.

Step 11: Define the SLA Required for Each Group

Primary support SLAs should include critical business units and functional group representation, such as networking operations, server operations, and application support groups. These groups should be recognized based on business needs as well as their part in the support process. Having representation from many groups also helps create an equitable overall support solution without individual group preference or priority. This can lead a support organization into providing premier service to individual groups, a scenario that may undermine the overall service culture of the organization. For example, a customer might insist his application is the most critical within the corporation when in reality the cost of downtime for that application is significantly less than others in terms of lost revenue, lost productivity, and lost customer goodwill.

Different business units within the organization will have different requirements. One goal of the network SLA should be agreement on one overall format that accommodates different service levels. These requirements are generally availability, QoS, performance, and MTTR. In the network SLA, these variables

are handled by prioritizing business applications for potential QoS tuning, defining help–desk priorities for MTTR of different network–impacting issues, and developing a solution matrix that will help handle different availability and performance requirements. An example of a simple solution matrix for an enterprise manufacturing company may look something like the following table. You can add information on availability, QoS, and performance.

Business Unit	Applications	Cost of Downtime	Problem Priority when Down	Server/Network Requirement
Manufacturing	ERP	High	1	Highest redundancy
Customer Support	Customer care	High	1	Highest redundancy
Engineering	File server, ASIC design	Medium	2	LAN core redundancy
Marketing	File server	Medium	2	LAN core redundancy

Step 12: Choose the Format of the SLA

The format for the SLA can vary according to group wishes or organizational requirements. The following is a recommended example outline for the network SLA:

1. Purpose of agreement
 - ◆ Parties participating in agreement
 - ◆ Objectives and goals for agreement
2. Services provided and products supported
 - ◆ Help–desk service and call tracking
 - ◆ Problem severity definitions based on business impact for MTTR definitions
 - ◆ Business–critical service priorities for QoS definitions
 - ◆ Defined solution categories based on availability and performance requirements
 - ◆ Training requirements
 - ◆ Capacity planning requirements
 - ◆ Escalation requirements
 - ◆ Reporting
 - ◆ Network solutions provided
 - ◆ New solution requests
 - ◆ Unsupported products or applications
3. Business policies
 - ◆ Support during business hours
 - ◆ After–hour support definitions
 - ◆ Holiday coverage
 - ◆ Contact phone numbers
 - ◆ Workload forecasting
 - ◆ Grievance resolution

- ◆ Service entitlement criteria
- ◆ User and group security responsibilities
- 4. Problem management procedures
 - ◆ Call initiation (user and automated)
 - ◆ First-level response and call repair ratio
 - ◆ Call tracking and record keeping
 - ◆ Caller responsibilities
 - ◆ Problem diagnosis and call-closure requirements
 - ◆ Network management problem detection and service response
 - ◆ Problem resolution categories or definitions
 - ◆ Chronic problem handling
 - ◆ Critical problem/exception call handling
- 5. Service quality goals
 - ◆ Quality definitions
 - ◆ Measurement definitions
 - ◆ Quality goals
 - ◆ Mean time to initiate problem resolution by problem priority
 - ◆ Mean time to resolve problem by problem priority
 - ◆ Mean time to replace hardware by problem priority
 - ◆ Network availability and performance
 - ◆ Managing capacity
 - ◆ Managing growth
 - ◆ Quality reporting
- 6. Staffing and budgets
 - ◆ Staffing models
 - ◆ Operations budget
- 7. Agreement Maintenance
 - ◆ Conformance review schedule
 - ◆ Performance reporting and review
 - ◆ Reconciliation of report metrics
 - ◆ Periodic SLA updates
- 8. Approvals
- 9. Attachments and exhibits
 - ◆ Call-flow diagrams
 - ◆ Escalation matrix
 - ◆ Network solution matrix
 - ◆ Report examples

Step 13: Develop SLA Workgroups

The next step is identifying participants in the SLA working group, including a group leader. The workgroup can include users or managers from business units or functional groups or representatives from a geographic base. These individuals communicate SLA issues to their respective workgroups. Managers and decision-makers who can agree on key SLA elements should participate. These individuals may include both managerial and technical individuals who can help define technical issues related to the SLA and make IT-level decisions (i.e., help desk manager, server operations manager, application managers, and network operations manager).

The network SLA workgroup should also consist of broad application and business representation in order to obtain agreement on one network SLA that encompasses many applications and services. The workgroup should have the authority to rank business–critical processes and services for the network, as well as availability and performance requirements for individual services. This information will be used to create priorities for different business–impacting problem types, prioritize business–critical traffic on the network and create future standard networking solutions based on business requirements.

Step 14: Hold Workgroup Meetings and Draft the SLA

The workgroup should initially create a workgroup charter. The charter should express the goals, initiatives, and time frames for the SLA. Next the group should develop specific task plans and determine schedules and timetables for developing and implementing the SLA. The group should also develop the reporting process for measuring the support level against support criteria. The final step is creating the draft SLA agreement.

The networking SLA workgroup should initially meet once a week to develop the SLA. After the SLA has been created and approved, the group may meet monthly or even quarterly for SLA updates.

Step 15: Negotiate the SLA

The last step in creating the SLA is final negotiation and sign–off. This step includes:

- Reviewing the draft
- Negotiating the contents
- Editing and revising the document
- Obtaining final approval

This cycle of reviewing the draft, negotiating the contents, and making revisions may take multiple cycles before the final version is sent to management for approval.

From the network manager's perspective, it is important to negotiate achievable results that can be measured. Try to back up performance and availability agreements with those from other related organizations. This may include quality definitions, measurement definitions, and quality goals. Remember that added service is equivalent to extra expense. Make sure that user groups understand that additional levels of service will cost more and let them make the decision if it is a critical business requirement. You can easily perform a cost analysis on many aspects of the SLA such as hardware replacement time.

Step 16: Measure and Monitor SLA Conformance

Measuring SLA conformance and reporting results are important aspects of the SLA process that help to ensure long–term consistency and results. We generally recommend that any major component of an SLA be measurable and that a measurement methodology be put in place prior to SLA implementation. Then hold monthly meetings between user and support groups to review the measurements, identify problem root causes, and propose solutions to meet or exceed the service level requirement. This helps make the SLA process similar to any modern quality improvement program.

The following section provides additional detail on how management within an organization can evaluate its SLAs and its overall service level management.

Service Level Management Performance Indicators

Service Level management performance indicators provide a mechanism to monitor and improve service levels as a measure of success. This allows the organization to react faster to service problems and to more

easily understand issues that impact service or the cost of down time in its environment. Not measuring service level definitions also negates any positive proactive work done because the organization is forced into a reactive stance. Nobody will call saying the service is working great, but many users will call saying the service is not meeting their requirements.

Service Level management performance indicators are therefore a primary requirement for service level management because they provide the means to fully understand existing service levels and to make adjustments based on current issues. This is the basis for providing proactive support and making quality improvements. When the organization does root-cause analysis on the issues and makes quality improvements, this then may be the best methodology to improve availability, performance, and service quality available.

For example, consider the following real scenario. Company X was getting numerous user complaints that the network was frequently down for extended periods of time. By measuring availability, the company found the major problem to be a few WAN sites. Closer investigation of those sites revealed that most of the problems were at a few WAN sites. The root cause was found and the organization resolved the problem. The organization then set service level goals for availability and made agreements with user groups. Future measurements identified problems quickly because of non-conformance to the SLA. The networking group was then viewed as having higher professionalism, expertise, and an overall asset to the organization. The group effectively moved from reactive to proactive in nature and helped the bottom line of the company.

Unfortunately, most networking organizations today have limited service level definitions and no performance indicators. As a result, they spend most of their time reacting to user complaints or problems instead of proactively identifying the root cause and building a network service that meets business requirements.

Use the following SLA performance indicators to determine the success of the service level management process:

- Documented service level definition or SLA that includes availability, performance, reactive service response time, problem resolution goals, and problem escalation
- Performance indicator metrics, including availability, performance, service response time by priority, time to resolve by priority, and other measurable SLA parameters
- Monthly networking service level management meetings to review service level compliance and implement improvements

Documented Service Level Agreement or Service Level Definition

The first performance indicator is simply a document detailing the SLA or service level definition. The primary goals of the service level definition should be availability and performance because these are the primary user requirements.

Secondary goals are important because they help define how the availability or performance levels will be achieved. For instance, if the organization has aggressive availability and performance targets, it will be important to prevent problems from occurring and to fix problems quickly when they occur. The secondary goals help define the processes needed to achieve the desired availability and performance levels.

Reactive secondary goals include:

- Reactive service response time by call priority
- Problem resolution goals or MTTR
- Problem escalation procedures.

Proactive secondary goals include:

- Device–down or link–down detection
- Network error detection
- Capacity or performance problem detection.

The service level definition for primary goals, availability, and performance should include:

The goal

- How the goal will be measured
- Parties responsible for measuring availability and performance
- Parties responsible for availability and performance targets
- Non–conformance processes

If possible, we recommend that the parties responsible for measurement and the parties responsible for results be different to prevent a conflict of interest. From time to time, it you may also need to adjust availability numbers because of add/move/change errors, undetected errors, or availability measurement problems. The service level definition may also include a process for modifying results to help improve accuracy and to prevent improper adjustments. See the next section for methodologies to measure availability and performance.

The service level definition for reactive secondary goals defines how the organization will respond to network or IT–wide problems after they are identified, including:

- Problem priority definitions
- Reactive service response time by call priority
- Problem resolution goals, or MTTR
- Problem escalation procedures

In general, these goals define who will be responsible for problems any given time and to what extent those responsible should drop their current tasks to work on the defined problems. Like other service level definitions, the service level document should detail how the goals will be measured, parties responsible for measurement, and non–conformance processes.

The service definition for proactive secondary goals defines how the organization provides proactive support, including the identification of network down, link–down or device–down conditions, network error conditions, and network capacity thresholds. Set goals that promote proactive management because quality proactive management helps eliminate problems and helps fix problems faster. This is normally accomplished by setting a goal of how many proactive cases are created and resolved without user notification. Many organizations set up a flag in help desk software to identify proactive cases versus reactive cases for this purpose. The service level document should also contain information on how the goal will be measured, parties responsible for measurement, and non–conformance processes.

Performance Indicator Metrics

We always recommend that any defined service level goal be measurable, allowing the organization to measure service levels, identify root–cause service issues that are inhibiting the primary goal of availability and performance, and make improvements that are aimed at specific targets. Overall, metrics are simply a tool that allows network managers to manage service level consistency and to make improvements according to business requirements.

Unfortunately, many organizations do not collect availability, performance, and other metrics. Organizations attribute this to the inability to provide complete accuracy, cost, network overhead, and available resources. These factors can impact the ability to measure service levels, but the organization should focus on the overall goals to manage and improve service levels. Many organizations have been able to create low-cost, low-overhead metrics that may not provide complete accuracy, but do satisfy these primary goals.

Measuring availability and performance is one area often neglected in service level metrics. Organizations that are successful with these metrics use two fairly simple methods. One method is to send Internet Control Message Protocol (ICMP) ping packets from a core location in the network to edges. You can also obtain performance using this method. Organizations that are successful with this method also group like devices into "availability groups," such as LAN devices or domestic field offices. This is also attractive because organizations usually have different service level goals for different geographic or business-critical areas of the network. This allows the metrics group to average all devices with the availability group to obtain a reasonable result.

The other successful method of calculating availability is to use trouble tickets and a measurement called impacted user minutes (IUM). This method tabulates the number of users that have been affected by an outage and multiplies it by the number of minutes of the outage. When expressed as a percentage of total minutes in the time period, this can be easily converted to availability. In either case, it can also be helpful to identify and measure the root cause of down time so that improvement can be more easily targeted. Root-cause categories include hardware problems, software problems, link or carrier problems, power or environment problems, change failures, and user error.

Measurable reactive support goals include:

- Reactive service response time by call priority
- Problem resolution goals, or MTTR
- Problem escalation time

Measure reactive support goals by generating reports from help desk databases, including the following fields:

- The time a call was initially reported (or entered into the database)
- The time the call was accepted by an individual working on the problem
- The time the problem was escalated
- The time the problem was closed

These metrics may require management influence to consistently enter problems in the database and update problems in real time. In some cases, organizations are able to automatically generate trouble tickets for network events or e-mail requests. This helps provide accuracy for identifying the start time of a problem. Reports generated from this kind of metric will normally sort problems by priority, work group, and individual to help determine potential issues.

Measuring proactive support processes is more difficult because it requires you to monitor proactive work and calculate some measurement of its effectiveness. Little work has been done in this area. It is clear, however, that only a small percentage of people will actually report network problems to a help desk, and when they do report the problem, it will clearly take time to explain the problem or isolate the problem as being network-related. Not all proactive cases will have an immediate effect on availability and performance either because of failure of redundant devices or links will have little impact on end users.

Organizations that implement proactive service level definitions or agreements do so because of business requirements and potential availability risk. Measurement is then done in terms of the quantity or percentage of proactive cases, as opposed to reactive cases that are generated by users. It is a good idea to measure the

amount of proactive cases in each area as well. These categories would include down devices, down links, network errors, and capacity violations. Some work may also be done using availability modeling and the proactive cases to determine the effect in availability achieved by implementing proactive service definitions.

Service Level Management Review

Another measure of service level management success is the service level management review. This should be done whether or not SLAs are in place. Perform the service level management review in a monthly meeting with individuals responsible for measuring and providing defined service levels. User groups may also be present when SLAs are involved. The purpose of the meeting is to then review performance of the measured service level definitions and to make improvements.

Each meeting should have a defined agenda that includes:

- Review of measured service levels for the given period
- Review of improvement initiatives defined for individual areas
- Current service level metrics
- A discussion of what improvements are needed based on the current set of metrics.

Over time, the organization may also trend service level compliance to determine the effectiveness of the group. This process is not unlike a quality circle or quality improvement process. The meeting helps target individual problems and determine solutions based on root cause.

Service Level Management Summary

In summary, service level management allows an organization to move from a reactive support model to a proactive support model where network availability and performance levels are determined by business requirements, not by the latest set of problems. The process helps create an environment of continuous service level improvement and increased business competitiveness. Service Level management is also the most important management component for proactive network management. For this reason, service level management is highly recommended in any network planning and design phase and should start with any newly defined network architecture. This allows the organization to implement solutions correctly the first time, with the least amount of downtime or rework.

Related Information

- **Technical Support – Cisco Systems**

All contents are Copyright © 1992–2005 Cisco Systems, Inc. All rights reserved. Important Notices and Privacy Statement.

Updated: Oct 04, 2005

Document ID: 15117
