# "Network performance optimization"

# بهینه سازی کارآیی شبکه



By

## Attaulhaq Bawar

**MIHE-2014-0549**

## Bachelor of Computer Science

*Under the Supervision of*

## *Mr.Mujtaba Jawed*

Assistant Professor

# Maiwand Institute of Higher Education

# Taimani Project, Kabul Afghanistan

Spring, 2017

Maiwand Institute of Higher Education

# "Network Performance Optimization"

A Thesis Presented to

Maiwand Institute of Higher Education

In partial fulfillment of the requirement for the degree of

# Bachelor in Computer Science

By

Attaulhaq Bawar

MIHE-2014-549

Spring 2017

# Final Approval

The undersigned have examined the thesis entitled '(Network Performance Optimization) presented by Attulhaq Bawar, a candidate for the degree of BCS (Bachelor of Computer Science) and hereby certify that it is worthy of acceptance.

_____                    _____

Signature & Date                                    Examiner

_____                    _____

Signature &Date                                     Supervisor

_____                    _____

Signature & Date                                    Co-Supervisor (Where Required)

_____                    _____

Signature & Date                                    HoD

_____                    _____

Signature & Date                                    Dean of Faculty

# DECLARATION

I Attaulhaq Bawar hereby declare that I have produced the network presented in this thesis, during the scheduled period of study. I also declare that I have no taken any material from any source except referred to wherever due that amount of plagiarism is within acceptable range. If a violation of university rules on research has occurred in this project, I shall be liable to punishable action under the plagiarism rules of the university.

 

 

**Attaulhaq Bawar**
**MIHE-2014-549**

# CERTIFICATE

It is certified that Attaulhaq Bawar has carried out all the work related to this thesis under my supervision at the Computer Science Faculty, Maiwand Institute and the work fulfills the requirement for award of Bachelor degree.

**Mr M. Jamshid Mahboob**
**HoD**

# DEDICATION

This Thesis is dedicated to:

My beloved parents and my great best friends who supported me encouraged me who leads me through the valley of darkness with lights of hope and support, inspiration and never stop giving of themselves in countless way stands by me during this precious period of my life.

# ACKNOWLEDGEMENTS

# ABSRACT

Network Optimization is the technology used for improving the performance for the Network. And it can be considered a precious component of effective information systems management as the technologies and techniques that are combined to gear towards improving network performance. It is necessary to use new standards and new devices for best performance in the network and avoid wastage of the time, data, and money for the business and saving benefits for a company. So the optimization can keep our network at the peak of efficiency that improves our network.

The goal of any Network optimization is to ensure a steady network with lowest cost, uptime data, and high availability. Optimized network must have the ability to ensure the best usage of the resources in that network. And also improve the availability for the customers to provide good services to them. So for network optimization make use of fault tolerance, redundancy, load balancing, reducing latency, congestion control, error control, flow control and traffic shaping for best transmission.

# TABLE OF CONTENTS

# LIST OF FIGURES

# CHAPTER I

# INTRODUCTION

Network optimization is technology or the way that is used for improving network performance for the given environment. And it is considered the best important component of effective information systems management. Network optimization plays an important role for the business in a company. New standards and devices making networks go faster than ever. Learn the best way to take advantage of them, saving your company time and money. We must optimize network performance. Network monitoring allows you to optimize your network performance, setting yourself apart in a highly competitive industry, and benefiting your team, partners, customers and contributing to the long terms success of your organization. Let's take a deeper look some reasons why keeping your environment operating at peak efficiency improves your business

The goal of any network optimization is to ensure an optimal network design with lowest cost structure and free flow of data. The optimized network should be able to ensure the usage of system resources, and also improve the availability for the customers. Network optimization makes use of traffic shaping, redundant data elimination, data caching and data compression and streaming of data protocols. Network optimization must be able to boost its efficiency without need of additional or price full hardware or the software.

There are many benefits of network optimization. It can help in faster data transfers including bulk data transfer, disaster recovery capabilities, reducing bandwidth expenses and also improving response times for interactive applications like databases and software applications. It also improves the performance of applications with better bandwidth and helps in maximizing network speeds between remote locations.

## 1.1 Fault tolerance

Fault tolerance it is the property that enables a system to continue the operating properly in the event of the failure of (or one or more faults within the system) some of its components. If the operation decreases or the quality decreases at all then the decrease is proportional to the severity of the failure, as compared to a naively designed

system in which even a small failure can cause total breakdown. Fault tolerance is particularly sought after in high-availability or life-critical systems. The ability of maintaining functionality when portions of a system break down if we have the fault tolerant gracefully we will be in success.

A design of fault tolerant enables a system to continue its intended operation, possibly at a reduced level rather than failing completely when some part of the system fails. The term is most commonly used to describe computer systems designed to continue more or less fully operational with perhaps, a reduction in through put or an increase in response time in the event of some partial failure. That is the system as a whole is not stooped due to problems either in the hardware or the software. An example in another field is a motor vehicle designed so it will continue to be drivable if one of the tires is punctured, or a structure that is able to retain its integrity in the presence of damage due to causes such as fatigue, corrosion, manufacturing flaws, or impact.

Within the scope of an individual system, fault tolerance can be achieved by anticipating exceptional conditions and building the system to cope with them and, and in general, aiming for self-stabilization so that the system converges towards an error-free state. Hoe ever if the consequences of a system failure are catastrophic, or the cost of making it sufficiently reliable is very high, a better solution may be to use some form of duplication. In any case, if the consequence of a system failure is so catastrophic, the system must be able to use reversion to fall back to a safe mode. This is similar to roll-back recovery but can be human action if humans are present in the loop.

The objective of creating a fault-tolerant system is to prevent disruptions arising from a single point of failure, ensuring the high availability and business continuity of mission-critical applications or systems.

Fault-tolerant systems use backup components that automatically take the place of failed components, ensuring no loss of service. These include figure 1

**Figure 1 - Fault-tolerant system backup**

## 1.2 Hardware systems Fault Tolerance

Most real time system must function with high availability even under hardware fault conditions. So real time systems are equipped with redundant hardware modules. Whenever a fault is encountered the redundant modules take over the functions of failed hardware module. Hardware fault tolerance is the most nature area in the general field of fault tolerant computing. Many hardware fault tolerance techniques have been developed and used in practice of critical applications ranging from telephone exchanges to space mission. In the past, the main obstacle to wide use of hardware fault tolerance has been the cost of extra hardware required. With the continued reduction in the cost of hardware, this is no longer a significant drawback, and notably on power consumption, may continue to restrict the use of massive redundancy in many applications. Hardware systems are backed up by identical or equivalent systems. For example, a server can be made fault tolerant by using an identical server running in parallel, with all operations mirrored to the backup server.

Hardware redundancy may be provided in one of the three following ways: one for one redundancy, N+X redundancy, load sharing.

**One for One Redundancy**

Here each hardware module has a redundant hardware module. The hardware module that performs under normal condition is called active and the redundant unit is called standby. The standby keeps monitoring the active unit at all times. It will take over and become active if the active unit fails. Since standby has to take over under fault conditions it has to keep itself synchronized with the active unit operations. Since the probability of both the units failing at the same time is very low, this technique

provides the highest level of availability. The main disadvantage here is that it doubles the hardware cost.

**N+X Redundancy**

In here if N hardware modules are required to perform system functions, the system is configured with N+X hardware modules; typically X is much smaller than N. whenever any of the N modules fails, one of the X modules takes over its functions. Since health monitoring of N units by X units at all times is not practical, a higher level module monitors the health of N units. If one of the N units fails, it selects one of the X units (it may be noted that one for one is a special case of N+X). The advantage lies in reduced hardware cost of the system as only X units are required to backup N units. However, in case of multiple failures, this scheme provides lesser system availability.

**Load sharing**

In here, under zero fault conditions, all the hardware modules that are equipped to perform system functions, share the load. A higher level module performs the load distribution. It also maintains the health status of the hardware units. If one of the load sharing modules fails, the higher level module starts distributing the load among the rest of the units. There is graceful degradation of performance with hardware failure. In here there is almost no extra hardware cost to provide the redundancy. The main disadvantage is that if a hardware failure happens during the busy hour, system will perform at a sub-optimal level until the failed module is replace. And network load balancing is a different flavor of load sharing where there is no higher level processor to perform load distribution. Instead the load distribution is achieved by hashing on the source address bits.

## 1.3 Software systems fault tolerance

Software fault tolerance is the ability for software to detect and recover from a fault that Is happening or has already happened in either the software or hardware in the system in which the software is running in order to provide service in accordance with the specification. Software fault tolerance is a necessary component in order to construct the next generation of highly available and reliable computing systems from embedded systems to data warehouse systems. So software fault tolerance is not a solution unto itself however, and it is important to realize that software fault tolerance

is just one piece necessary to create the next generation of systems. In order to understand software fault tolerance it is important to understand the nature of the problem that software fault tolerance is supposed to solve. Software faults are all design faults. The source of the problem being solely design faults is very different than almost any other system in which fault tolerance is desired property. This inherent issue that software faults are the result of human error in interpreting a specification or correctly implementing an algorithm, creates issue which must be dealt with in the fundamental approach to software fault tolerance. Currently software fault tolerance methods are based on traditional hardware fault tolerance. The deficiency with this approach is that traditional hardware fault tolerance was designed to conquer manufacturing faults primarily, and environmental and other faults secondarily. Design diversity ws not a concept applied to to the solutions to hardware fault tolerance. Software tolerance tries to leverage the experience of hardware fault tolerance to solve different problem, but by doing so creates a need for design diversity in order to properly create a redundant system. Software systems are backed up by other software instances. For example, a database with customer information can be continuously replicated to another machine. If the primary database goes down, operations can be automatically redirected to the second database.

## 1.4    Power sources fault tolerance

Power sources that are made fault tolerant using alternative sources. For example, many organizations have power generators that can take over in case main line electricity fails. So equipment used in critical applications must contain a power supply that is fault tolerant in order to guarantee continuous and uninterrupted functionality. We will focus on the importance of system power supply. It is widely accepted that the highest stresses in any electronic system is mostly in its power supply. This is where most of the power losses exist and therefore the highest thermal stress. The input derived from the utility (or a generator) is subject to fluctuations, brownouts, transients and interruptions. The power components in the power supply are switched at high speed while conducting high current or subjected to high voltages. Further the switching speeds are ever increasing in order to reduce volume and weight. All of these factors make the power supply of any electronic system the most vulnerable to failure, resulting in total system dysfunction. This is therefore the main reason for designing the system with backup features which permit uninterrupted power delivery despite one

more failures in the utility or power processing path. In the ultimate case an FTPS is a power supply ( with DC or AC input or output) which delivers uninterrupted output to equipment (or functional circuits) despite one or more of the following detrimental events:

1.      A failure of the input power source (utility or generator)

2.      A failure of the power supply itself

3.      A failure of the battery used as a back-up energy source within the backup apparatus.

In similar fashion, any system or component which is a single point of failure can be made fault tolerant using redundancy.

Fault tolerance can play a role in a disaster recovery strategy. For example, fault-tolerant systems with backup components in the cloud can restore mission-critical systems quickly, even if a natural or human-induced disaster.

## 1.5     Fault Tolerance vs. High Availability

High availability refers to a system's ability to avoid loss of service by minimizing downtime. It's expressed in terms of a system's uptime, as a percentage of total running time. Five nines, or 99.999% uptime, is considered the "holy grail" of availability.

In most cases, a business continuity strategy will include both high availability and fault tolerance to ensure your organization maintains essential functions during minor failures, and in the event of a disaster.

While both fault tolerance and high availability refer to a system's functionality over time, there are differences that highlight their individual importance in your business continuity planning.

Consider the following analogy to better understand the difference between fault tolerance and high availability. A twin-engine airplane is a fault tolerant system – if one engine fails, the other one kicks in, allowing the plane to continue flying. Conversely, a car with a spare tire is highly available. A flat tire will cause the car to stop, but downtime is minimal because the tire can be easily replaced.

Some important considerations when creating fault tolerant and high availability systems in an organizational setting include:

## 1.6    Downtime

A highly available system has a minimal allowed level of service interruption. For example, a system with "five nines" availability is down for approximately 5 minutes per year. A fault-tolerant system is expected to work continuously with no acceptable service interruption.

## 1.7    Scope

High availability builds on a shared set of resources that are used jointly to manage failures and minimize downtime. Fault tolerance relies on power supply backups, as well as hardware or software that can detect failures and instantly switch to redundant components.

## 1.8    Cost

A fault tolerant system can be costly, as it requires the continuous operation and maintenance of additional, redundant components. High availability typically comes as part of an overall package through a service provider (e.g., load balancer provider).

Some of your systems may require a fault-tolerant design, while high availability might suffice for others. You should weigh each system's tolerance to service interruptions, the cost of such interruptions, existing SLA agreements with service providers and customers, as well as the cost and complexity of implementing full fault tolerance.

## 1.9    Redundancy

Redundancy is one of the layman's terms backup and it is the right option to be chosen for a network design and its availability as it is the best technique for the fault tolerant mechanism. There are also several different types of redundancy: Network, hardware, power, and geographic. As we know that web hosting is the service that it allows organizations and individuals to post a website or web page on to internet and it can be for the reason of business and company will have multiple layers of redundancy to ensure that their data is safe and to maximize their service uptime.

# CHAPTER II

## LATENCY AND DELAY IN RESPONSE TIME

Latency is a networking term that describe the total time it takes a data packet to travel from one node to another. In other word when the data packet is transmitted and returned back to its source, the total time for a round trip is known as latency. Latency refers to any of several kinds of delays typically incurred during the data processing in the network. So and it can be refers to the time interval when a system components is waiting for the another components to do something this duration is called latency. The low latency of the network connection is one of that experiences small delay times, while a high latency connection suffers from long delays. In data communication, digital networking and packet-switched networks, latency is used in two major contexts. One represents a one-way trip while the other is a round trip. One-way latency is measured by counting the total time it takes a packet to travel from its source to its destination. Round-trip latency is measured by adding one-way latency from the destination to the time it takes the packet to return from the destination and arrive back at the source. Unlike one-way latency, round-trip latency always excludes processing time at the destination point. A service called ping is used to measure round-trip latency. In formal network transmission, the following four elements are involved in latency.

## 2.1    Delays in Storage

As data is written on hard disks and other storage devices, a delay occurs in reading and writing to and from different blocks of memory. Processors often consume a lot of time finding the exact location for reading and writing data. Sometimes intermediate devices like switches or hubs also cause delays.

## 2.2    Device Process

Latency is not limited to storage devices but can also be caused by different network devices. For example, when a router receives a data packet, it keeps that packet for a few seconds to read its information and also to write some extra information.

## 2.3    Transmission

There are many kinds of transmission media and all have limitations. Each medium, from fiber optics to coaxial cables, takes some time to transmit one packet from a source

to a destination. Transmission delays depend on packet size; smaller packets will take less time to reach their destination than larger packets.

## 2.4    Propagation

Delays occur even when packets travel from one node to another at the speed of light. So propagation delays, latency may also involve transmission delays(properties of the physical medium) and processing delays ( such as passing through the servers or making network hops on the internet).

Latency in communication can be created in live transmissions from various points on the earth as the communication hops between a ground transmitter and a satellite and from a satellite to a receiver each take time. Consumers of networks connecting from one distance to another these live events can be seen to have to wait for responses. The latency can be the wait time introduced by the signal traveling all the geographical distance as well as over the various pieces of communications equipment. Even fiber optics is limited by more than just the speed of light, as the refractive index of the cable and all repeaters or amplifiers along their length introduce delays

## 2.5    Latency types

Network latency is the word or expression of time that means how much it takes for a packet of the data to reach from one designated point to another. In some environments as example, in (AT&T), latency is measured by sending a data packet that is returned to the sender; then the round trip time is considered the latency. Ideally, latency is close to zero as possible.

## 2.6    Internet latency

It is just a special case of network latency – we know that internet is a very large wide area networks (WAN). The same factors as above determine latency on the internet. However the distance in the transmission medium, the number of hops over equipment and servers are all greater than smaller networks. The internet latency measurement would generally start at the exit of a network and end on the return of the requested data from an internet resource.

## 2.7    Interrupt latency

It is the length of the time that it takes for computer to act on an interrupt, which is a signal telling the operating system to stop until it can decide what it should do in response to some event.

## 2.8    WAN Latency

It can be an important factor in determining the internet latency. A Wan that can be busy for directing other traffic will produce a delay whether a resource is being requested from the server on the LAN, other computers on that network or elsewhere on internet. LAN users will also experience delay when the WAN is busy. In either of these examples the delay would still exist if the rest of the hops – including the server where the desired data was located- were entirely free of traffic congestion.

## 2.9    Audio latency

It is the delay between sound being created and heard. As the sound created in the physical world, this delay is determined by the speed of sound, which varies depending on the medium the sound wave travels through. The sound can be travels faster in denser mediums: It travels faster through solids, less quickly through liquids and slowest through air. We generally refer to the speed of sound that measured in dry air at room temperature, which is 796 miles per-hour. In the electronics, audio latency is cumulative delay from input audio to output audio. And this delay depends on the hardware and even software used, such as the operating system and drivers used in computer audio. The latencies of the 30 milliseconds are generally noticed by an individual as a separate production and arrival of the sound to the ear.

## 2.10   Computer and operating system latency

It is the combined delay between an input or command and the desired output. In the computer system, latency is often used to mean any delay or waiting that increases real or perceived response time beyond what is desired. Specific contributors to computer latency.

## 2.11   Latency testing

Latency testing can vary from application to application. In some of applications, measuring latency requires special and complex equipment or knowledge of special computer commands and programs; in other cases, latency can be measured with a stop

watch. In networking, an estimated latency to equipment or servers can be determined by running a ping command; information about latency through all the hops can be gathered with a trace route command. High speed cameras might be used to capture the minute differences in response times for input to various mechanical and electronic systems.

## 2.12   How to Reduce Latency

Ideally, data centers and disaster recovery sites should be placed at a longer distance from each other than is traditionally practiced to literally insure your data from the impact of any man-made or natural disaster. As Data is the lifeblood of business specially in networking environment. So a slow data transfer rate makes it harder to analyze, back-up, and restore data. Many organizations have to battle data latency on a daily basis, hampering their ability to deliver new digital products and services, be profitable, handle customer relationships, and retain operational efficiency. Data latency is a serious business issue that needs to be addressed. In contrast, network latency is a technical issue; but they both correlate with each other.

## 2.13   Tackling Latency

There is very little you can do to reduce network latency. The only way you can reduce it is by moving data centers or disaster recovery sites close to each other. This has often been the traditional approach, but from a disaster recovery perspective, it can be disastrous because each of your data centers could end up being located in the same circle of disruption. Ideally, data centers and disaster recovery sites should be placed at a longer distance from each other than is traditionally practiced to literally insure your data from the impact of any man-made or natural disaster.

Yet companies have to move data further and further away at ever increasing network speeds. The latency within the data center is very small these days, and so latency has its greatest impact on data transfer rates when it is moved outside of the data center to the cloud, or the internet for the benefit of customers.

The response by many organizations is to address latency issues with the implementation of traditional WAN optimization tools, which have little impact on latency and data acceleration. Another strategy to reduce latency is to increase the organization's bandwidth with a high capacity pipe, but again, this won't necessarily

accelerate the data, reduce latency and packet loss to the required levels. Yet conversely, it is only possible to mitigate the effects of latency to accelerate data over a WAN when data is being moved over long distances.

## 2.14 Traditional Response

Traditional WAN optimization vendors give the impression of reduced latency by keeping a copy of the data locally, so the perception is that latency has been reduced because you aren't going outside the data center. However, real latency still exists because you have two points going out across a WAN. More to the point, things have changed over the last 10 years. Eighty percent of data used to be generated and consumed internally and 20 percent externally. Disasters tended to be localized, and organizations had to cope with low bandwidth, low network availability, and a high cost per megabit. The data types were also highly compressible, involving small data sets. WAN optimization was, therefore, the solution because it used a local cache, compression and duplication, and locally optimized protocols.

## 2.15 Changing Trends

Today, everyone has moved from slow speed connections where everything was compressed, to using big pipes now that the price has come down. In the past, data was produced faster than the pipe could manage data flows, but now we can accelerate big volumes of data without having local caches. This is why many companies are transitioning from WAN optimization to WAN acceleration.

It should also be noted that the data scenario of 10 years ago has been reversed. Only 20 percent of data is now generated and consumed internally because 80 percent of it now emanates from external sources. Disaster often causes a wider impact today, and organizations have to cope with ever-increasing data sets, which are created by the growing volumes of big data. Another recent trend is the increased use of video for videoconferencing, marketing and advertising purposes.

Firms are also now enjoying higher bandwidth, increased availability, and lower cost per Mb networks. Files tend to be compressed, duplicated and encrypted across globally dispersed sites. This means that a new approach is needed that doesn't require a local cache, offers no data change, and uses any protocol to accelerate and reduce both data and network latency.

## 2.16   Strong Encryption

Converged systems can nevertheless help to address these issues, but strong encryption is needed to shorten the window of opportunity for hackers to intercept data flows. So if you are being attacked, you can quickly move data offsite. This can be achieved with tools that use machine intelligence to accelerate data that needs to travel across long distances without being changed.

With data acceleration and mitigated latency, it becomes possible to situate data centers and disaster recovery away from each other and outside their own circles of disruption. This approach offers a higher degree of security because large volumes of data can be transferred within minutes and seconds, denying a hacker any chance to do serious harm.

Organizations should, therefore, look beyond the traditional WAN optimization players to smaller and more innovative ones that can mitigate data and network latency whenever data is to be transmitted, received and analyzed over long distances. The future competitiveness, profitability, customer relationships and efficiency may depend on it. So it's time to look anew at latency to accelerate your data no matter where it resides – in the data center or outside of its walls.

## 2.17   Network Redundancy

Network redundancy is the process through which additional or the alternate instances of the network devices, equipment, protocols and communication mediums are installed within the network infrastructure. Redundancy protocols are the standards for that to provide seamless failover against failure of any network component. These redundancies are invisible to the application. And these protocols are installed in tow nodes of ports and are attached to two separated networks of similar topology. This method is for the ensuring network for a high availability in case of avoiding the network devices or path failure and unavailability. Network redundancy is primarily implemented in enterprise networks infrastructure to provide a redundant source of network communications. It will serves as a backup mechanism for quickly swapping network

Operation on the redundant infrastructure in the event of occurring unplanned network outages. Typically, network link Redundancy or Network Redundancy is achieved by

the addition of alternate networks paths, which are implemented through redundant standby routers and switches. When the primary path is unavailable or breakdown the alternate path can be instantly deployed to ensure minimal downtime and continuity services. For example a backup is simply for the network a web hosting company uses to get you online. Reliable providers have multiple internet carriers that their servers use to ensure that your website never goes down. A first hop redundancy protocol (FHRP) is a computer networking protocol which is designed to protect the default gateway used on sub network by allowing two or more routers to provide backup for that address, in the event of failure of an active router, the backup router will take the address, usually within a few seconds. Such protocols can also be used to protect other services operating on single IP address, not just a router

## 2.18  HSRP

Hot standby router protocol (HSRP) is a cisco proprietary it is the redundancy protocol for establishing a fault tolerant default gateway. Its version 1 describes in RFC 2281. And there is no RFC for version 2 of the protocol. Version 2 of such protocol introduces stability, scalability and diagnostic improvements. It is not compatible with version 1 HSRP. There is no RFC for version 2 of the protocol.

HSRP will establish the framework between the routers in order to achieve default gate way failover if the primary gateway becomes inaccessible or goes down. HSRP routers send multicast hello messages to other routers to notify them of their priorities (which router is preferred) and the current status (Active or stand by).

The primary router with the highest configured priority will act as a virtual router with a pre-defined gateway IP address and will respond to the ARP / ND request from machines connected to the LAN with a virtual MAC address. If the primary router should fail, the router with the next-highest priority would take over the gateway IP address and answer ARP requests with the same MAC address, thus achieving transparent default gateway failover.

HSRP is not a routing protocol as it does not advertise IP routes or affect the routing table in any way.

HSRP has the ability to trigger a failover if one or more interfaces on the router go down. This can be useful for dual branch routers each with a single link back to the

gateway. If the link of the primary router goes down, the backup router will take over the primary functionality and thus retain connectivity to the gateway. Figure 3



Figure 2- HSRP

So in HSRP we have one Active Router, one Standby Router, other routers are in standby group. And communication protocols are IP v4 with port UDP 1985 and Ipv6 Port UDP port 2029. Authentication default no authentication plain text authentication MD5 authentication. Active selector one router is selected as active, another as standby router. The remaining are in listen state. Highest values wins default value is 100. Hello and hold time r Hello interval between successive HSRP hello messages from a given router: default is 3 sec and hold interval between receipt of a hello, and the presumption that sending router failed so default is 10 sec. active timer is 10 sec and standby timer is 10 sec no use of preemption by default is off.

## 2.19  VRRP

VRRP or virtual router redundancy protocol is a computer networking protocol that provides an automatic assignment of available Internet Protocol (IP) routers to participating hosts. This increases the availability and reliability of routing paths via automatic default gateway selections on an IP sub network.

The protocol achieves this by creation of virtual routers, which are an abstract representation of multiple routers, i.e. master and backup routers, acting as a group. The default gateway of a participating host is assigned to the virtual router instead of a

physical router. If the physical router that is routing packets on behalf of the virtual router fails, another physical router is selected to automatically replace it. The physical router that is forwarding packets at any given time is called the master router.

VRRP provides information on the state of a router, not the routes processed and exchanged by that router. Each VRRP instance is limited, in scope, to a single subnet. It does not advertise IP routes beyond that subnet or affect the routing table in any way. VRRP can be used in Ethernet, MPLS and token ring networks with Internet Protocol Version 4 (IPv4), as well as IPv6.

The protocol is described in Internet Engineering Task Force (IETF) publication RFC 5798, which is an open standard. Figure 4



**Figure 3- VRRP**

## 2.20  GLBP

A (GLBP) or the Gateway Load Balancing Protocol it is a Cisco proprietary protocol that attempts to overcome the limitations of existing redundant router protocols by adding basic load balancing functionality.

In addition to being able to set priorities on different gateway routers, GLBP allows a weighting parameter to be set. Based on this weighting (compared to others in the same virtual router group), ARP requests will be answered with MAC addresses pointing to

different routers. Thus, by default, load balancing is not based on traffic load, but rather on the number of hosts that will use each gateway router.

GLBP elects one AVG or (Active Virtual Gateway) for each group. Other group members act as backup in case of AVG failure. In case there are more than two members, the second best AVG will place in the Standby state and all other members are placed in the Listening state. This is monitored using hello and hold time timers, which are 3 and 10 seconds by default. The elected AVG then assigns a virtual MAC address to each member of the GLBP group, including it, thus enabling AVFs (Active Virtual Forwarders). Each AVF assumes responsibility for forwarding packets sent to its virtual MAC address. There could be up to four AVFs at the same time.

By default, GLBP routers use the local multicast address 224.0.0.102 to send hello packets to their peers every 3 seconds over UDP 3222 (source and destination).



**Figure 4- GLBP**

## 2.21 CARP

The (CARP) Common Address Redundancy Protocol is a computer networking protocol which allows multiple hosts on the same range of local area network to share a set of IP addresses. Its purpose is to provide the failover redundancy, especially when used with firewalls and routers. In some configurations, the CARP can also provide load balancing functionality. CARP provides functionality similar to VRRP and to the Cisco Systems' HSRP. It is implemented in several BSD-based operating systems and has been ported

to the Linux.

Example for the CARP

If there is a single computer running a packet filter, and it goes down or become inactive, the networks on either side of the packet filter can no longer communicate with each other, or they communicate without any packet filtering. If, there are two computers running a packet filter, running CARP, then if one fails, the other will take over, and computers on either side of the packet filter will not be aware of the failure, so operation will continue as normal. In order to make sure the new master operates the same as the old one, the packet filter used must support synchronization of state between the two computers.

## 2.22  NSRP

NSRP or the Net Screen Redundancy protocol it is the proprietary of the juniper Networks router redundancy protocol providing the load balancing for the networks

To function properly as a network firewall, a security device must be placed at the single point through which all inter zone traffic must pass. When a single device is responsible for handling all inter zone traffic, it becomes vital that the traffic flow remain uninterrupted, even in the event of a device or network failure.

To ensure a continuous traffic flow, you can cable and configure two security devices in a redundant cluster, where one device acts as the primary and the other as its backup. The primary device propagates all its network and configuration settings, along with the current session information, to the backup. If the primary device fails, the backup becomes the primary and begins processing traffic.

You can configure a redundant cluster in one of three ways:

## 2.23  ESRP

Stands for the Extreme Standby Router Protocol it is the protocol created by Extreme Networks to provide redundancy for both layer 2 and layer3  as the VRRP used in conjunction with a layer 2 loop prevention protocol such as spanning Tree provides the functionality of ESRP.

## 2.24   R-SMLT

The full name of the RSMLT is Routed Split Multilink Trunking. It is link layer protocol which ensures that the traffic gets across a VLAN by providing redundant switches and redundant links whereas RSMLT ensures traffic can be routed off the VLAN by adding router redundancy. This protocol is developed at Nortel. The R-SMLT protocol works with SMLT and distributes Split Multi-Link Trunking (DSMLT) technologies to provide sub-second failover so no outage is noticed by end users. It also provide an active-active router concept to core SMLT networks

## 2.25   The importance of Redundancy

Today's networks are high-tech and most times high speed. Common to most wide area network (WAN) designs is the need for a backup to take over in case of any type of failure to your main link. Implementing a network redundancy strategy will depend on many factors, largely dictated by the application and existing network topology- the physical layout, location of systems, processes and devices, and the way the cabling infrastructure is run. Certain redundancy methodologies are more suited for system configuration than another. Redundancy protocols may be standards-based or proprietary. In general, standards-based redundancy protocols provide outstanding interoperability but slower recovery times, whereas proprietary protocols in most cases offer faster recovery speeds and are designed especially for industrial recovery applications. There is some overlap in the features and functionality of redundancy protocols, and in many applications, the use of hybrid protocols( a mix of methodologies) is quite common. Selecting the right protocol and configuration of redundancy protection requires careful evaluation of the application requirement and review of available redundancy options. A simple scenario would be if you had a single connection from your core site to each remote office or branch office you connect with. What if that link went down? How would you continue your operations if it did? In this section we will explore this scenario and other scenarios to help you design and plan for a backup solution that you can count on and one that is cost effective and will not break the bank.

Network redundancy is a simple concept to understand. If you have single point of failure and it fails you, then you have nothing to rely on. If you put in a secondary (or tertiary) method of access, then when the main connection goes down, you will have a

way to connect to resources and keep the business operational. The first step in creating network redundancy is to set up a project plan that will allow you to scrutinize the current architecture/infrastructure, plan for a way to make it redundant, plan for a way to deploy it and then setup a way to test it. Nothing should be thought of as complete until you have tested everything for operational success. Our final step will be putting in policy and process that allow you to monitor it and be alerted when things do fails so it can take action. Commonly a company's security policy, disaster recovery plan, business continuity plan and/or incident response plan will leave room for this type of solution.

Figure 1 shows the layout of common design, where a branch office needs to connect to a core site where centralized resources are located such as financial applications, enterprise resource planning (ERP) software, database, file server data and so on



**Figure 5- Network redundancy**

Figure 1: Shows a common WAN connection scenario with redundancy in place

Here you can see how the remote branch office connects to the core site. Here, there is a dedicated MPLS circuit/link that provides bandwidth at approximately 1.5 Mbps which is the connection speed of a T1. The MPLS routers are connected to the network via a network switch. Commonly, the router is the networks default gateway, where all packets are sent that are not found locally. The router needs to make a routing decision and since the mani link is up, decides to send via the MPLS connection. When the main

link drops, commonly the alternate link is used if set up correctly. This provides your remote site with new path to reach the resources needed to continue operations.

## 2.26  Load balancing

While Server Load Balancing (SLB) could mean many things it is refers to efficiently distributing incoming network traffic across a group of back end servers or server's pool or can be defined as a process and technology that distribute site traffic among several servers using a network-based device. This device intercepts traffic destined for a site and redirects that traffic to various serves. The load balancing process is completely transparent to the end user. There are often dozens or even hundreds of servers operating behind a single URL. In the figure 7.2, we see the simplest representation of SLB. Modern high-traffic websites must serve hundreds of thousands, if not millions, of concurrent requests from users or clients and return the correct text, images, video, or application data, all in a fast and reliable manner. To cost-effectively scale to meet these high volumes, modern computing best practice generally requires adding more servers.



**Figure 6- Server load balancing**

A load balancer acts as the "traffic cop" sitting in front of your servers and routing client requests across all servers capable of fulfilling those requests in a manner that maximizes speed and capacity utilization and ensures that no one server is overworked, which could degrade performance. If a single server goes down, the load balancer redirects traffic to the remaining online servers. When a new server is added to the server group, the load balancer automatically starts to send requests to it.

## 2.27 Load balancer functions

1.     Based on the chosen method of distributing network/internet traffic, the load balancer will deliver requests to the best network servers as fast and efficiently as possible

2.     Continually check the performance of the network servers and make decisions which server is performing in the best way to serve the users demands.

3.     Distributes client requests or network load efficiently across multiple servers

4.     Ensures high availability and reliability by sending requests only to servers that are online.

5.     Provides the flexibility to add or subtract servers as demand dictate.

## 2.28 Load Balancing Algorithms

Different load balancing algorithms provide different benefits; the choice of load balancing method depends on your needs: as we have three types of algorithms for load balancing .

1.     Round Robin

2.     Least Connection

3.     IP Hash

## 2.29 Round Robin

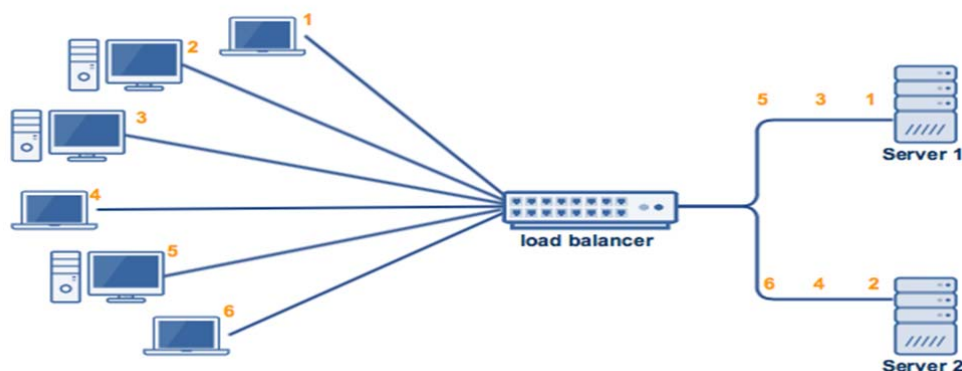Requests are distributed across the group of servers sequentially.



**Figure 7- Round Robin**

The most basic of load balancing services is called round robin. round robin load balancing allows you to distribute client requests across multiple servers. Load balancers improve server fault tolerance and end-user response time. As load balancing distribute client request s across multiple servers to optimize resources utilization. In example with a limited number of servers providing service to large number of clients, a server can become overloaded and degrade the server performance. Load balancing is used to prevent bottlenecks by forwarding the client requests to the server best suited to handle them.

In a load balancing setup, the load balancers are logically located between the client and the server farm. Load balancing is used to manage traffic flow to the servers in the server farm. Load balancing can be performed on HTTP,FTP, SSL_TCP,UDP,DNS, DNS-TCP, RSTP, and SSL-BRIDGE.

Load balancing uses a number of algorithms, called load balancing methods. To determine how distribute the load among the servers. When a load balancers are configured to use the round robin method. It rotates incoming requests around to the managed servers, regardless of the load.

When to use the round robin load balancing? When you we have decided that you need to deploy load balancers to improve the performance of our network servers we will also need to decide which method of server performance selection works best. The load balancers will need to be configured to use load balancing technique that takes the decision from one server in the cluster to decide where to direct traffic to dedicate load balancer that sits between the users and the servers and acts as kind of traffic cop monitoring server performance and directing traffic to the servers that offer the best response.

So deploying of load balancers allow us to use a load balancing technique that reduce greatly the risk and overload of traffic so means that our network will be optimized.

## 2.30   Round robin mechanism

It is the simple mechanism in which the content access request is responded to by the load balance in a rotational basis, the first request grants access to the first available content server giving its IP address and the second server IP address and so on. The moment a server IP address has been given its IP address is moved and to the back of

the list of available IP addresses and gradually it moves back to the top of the list and becomes available again. The frequency that it returns to the top depends on the number of available servers in the round robin server cluster being used. A good way to think of this is a method of server allocation on a continuous looping fashion. An example of round robin load balancing environment

Geographically distributed web servers are best served by Applying DNS load balancing round robin server content distribution. As an example a company c have a single domain name and four absolutely identical company home pages on four physical servers based in Herat, Mazar, nangarhar and badakhshan. Each time a request for access comes in the first one will be sent to the Herat server the second to the one in Mazar and then the third request goes to the Nangarhar server, as each IP address is given out it is automatically placed at the back of the list.

## 2.31  Least Connections

Performance is important, and that mean it is important that our infrastructure support the need for speed. Load balancing algorithms are an integral piece of performance equation and can both improve – or degrade- performance. That is why it is important to understand more about algorithms than their general selection mechanism. Understanding that round robin is basically an iterative choice, traversing a list one by one is good – but understanding what that means in term of performance and capacity on different types of application and its workloads is even better.

That is checked out "fastest response time" and today we are diving into "least connections" which, as stated above, does not mean "least loaded".

The industry standard "least connections" load balancing algorithm uses the number of current connections to each application instance (member) to make its load balancing decision. The member with the least connection is chosen.

The premise of this algorithm is a general assumption that fewer connections (and thus fewer users) mean less load and therefore better performance. That is operational at work – if performance decreases as load increase it stands to reason that performance increases as load decreases.

That can be true when all applications workloads required the same resources.

Unfortunately, that is no longer true and the result is uneven load distribution that leads to unpredictable performance fluctuations as demand increases. We consider a simple example: a user logging into system takes at least one if not more database queries to validate credentials and then update the sytem to indicate the activity. Depending on the nature of the application, other applications activities will require different quantities of resources. Some are RAM heavy Other CPU heavy, other files or database heavy. Furthermore, depending on the user in the question, the usage pattern will vary greatly. One hundred users can be logged into the same system (required at a minimum ten connections) but if they are all relatively idle, the system will be lightly loaded and performing well. Conversely, another instance may boast only 50 connections, but all fifty more heavily loaded and performance may be already beginning to suffer. When the next request comes in however, the load balancer using a "least connections" algorithm will choose the latter member, increasing the burden on that member and likely further degrading performance. The premise of the least connections algorithm is that the application instance with the fewest number of connections is the least loaded.

The only way to know which application instance is the least loaded is to monitor its system variables directly, gathering CPU utilization and memory and comparing it against known maximums. That generally requires either SNMP, agents, or other active monitoring mechanisms that can unduly tax the system in and of it by virtue of consuming resources. So we can say a least connection algorithm

A new request is sent to the server with the fewest current connections to clients. The relative computing capacity of each server is factored into determining which one has the least connections.
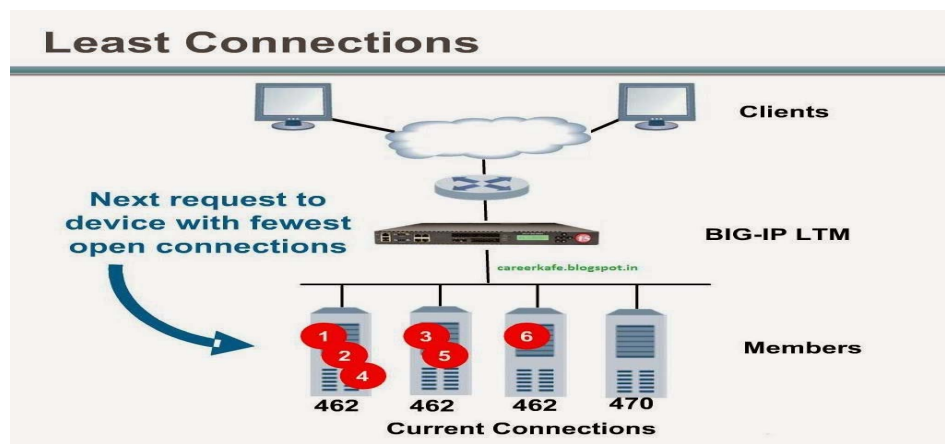


**Figure 8- Least connection**

## 2.32   IP Hash

IP hash load balancing uses an algorithm that takes the source and destination IP address of the client and server to generate a unique hash key. This key is used to allocate the client to particular server. As the key can be regenerated if the session is broken, this method of load balancing can ensure that the client is directed to the same server that it was using previously. This is useful if it is important that a client should connect to a session that active after a disconnection and reconnection. The IP address of the client is used to determine which server receives the request. So this algorithm is deterministic. It means that if no elements involved in the hash computation, then the result will be same 2 equipment are able to apply the same hash, load- balance the same way, making load- balancer failover transparent. A hash function is applied on the source IP address of the incoming request. The hash must take into account the number of servers and each server's weight. The main issue with source IP hash load balancing algorithm is that each change can redirect everybody to different server. That is why some good load-balancers have implemented a consistent hashing method which ensure that if  a server fails, for example only the client connected to this server are redirected. Consistent hashing is not providing perfect hashing, so in a farm of 4 servers, some may receive more clients than others. When a failed server comes back, its users will be redirected to it.
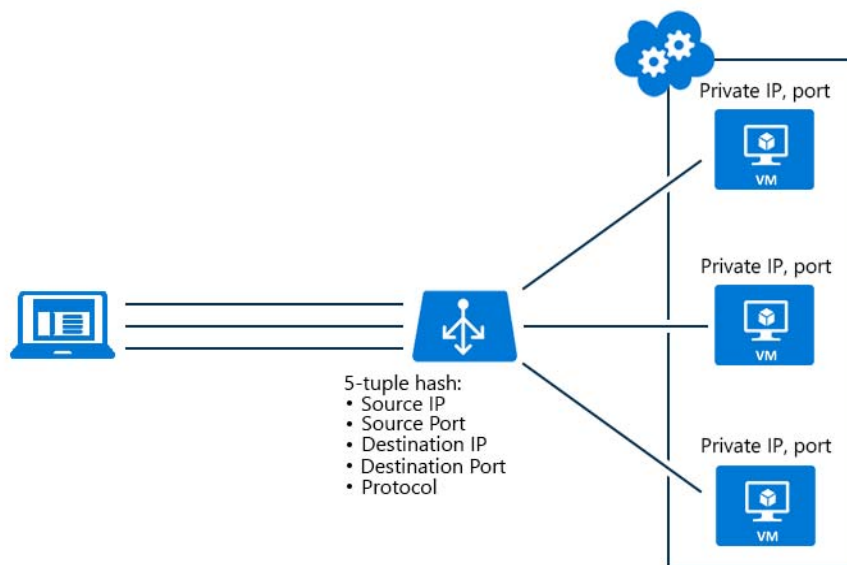


**Figure 9- IP Hash**

## 2.33   Error Detection and Correction

Network must be able to transfer data from one device to another with acceptable accuracy for most applications a system must guarantee that the data received are identical to the data transmitted. Any time data are transmitted from one node to the next the can become corrupted in passage. Many factors can alter one or more bits of a massage. Some application requires a mechanism for detecting and correcting errors. Some of application can tolerate a small level of errors for example random errors in audio or video transmissions may be tolerable, but when we transfer text we expect a very high level of accuracy.

Let me first discuss some issues related, directly or indirectly, to error detection and correction. Whenever bits flow from one point to another, they are subject to unpredictable changes because of interference. This interference can change the shape of the signal. In a single- bit error, a (0) is changed to a(1) or a(1) to a(0). In a burst error multiple bits are changed. For example a 1/100s burst of impulse noise on a transmission with a data rate of 1200 bps might change all or some of 12 bits of information.

## 2.34   Single-Bit Error

The term single –bit error means that only 1 bit of given data unit (such as a byte, character, or packet) Is changed from 1 to 0 or from 0 to 1. In a single-bit error, only 1 bit in the data unit has changed. To understand the impact of the change , imagine that each group of 8 bits is an ASCII character with a 0 bit added to the left, 00000010 (ASCII STX) was sent, meaning  start of text, but 00001010 (ASCII LF) was received, meaning line feed Single-bit errors are least likely type of error in serial data transmission. To understand why, imagine data sent data 1 Mbps. This means that each bit lasts only 1/1,000,000s, or 1u for a single-bit error to occur, the noise must have a duration of only 1us, which is very rare; noise normally lasts much longer than this.

## 2.35   Burst Error

The term burst error means that 2 or more bits in the data unit have changed from 1 to 0 or from 0 to 1. A Burst Error means that 2 or more bits in the data unit have changed. Shows the effect of a burst error on a data unit in this case, 0100010001000011 was sent, but 0101110101100011 were received. It is to mention that burst errors occur in consecutive bits. The length of burst is measured from the first corrupted bit to the last

corrupted bit. Some bits in between may not have been corrupted A burst error is more likely to occur than a single-bit error. The duration of a noise is normally longer than the duration of 1 bit which means that when noise affects data, it affects a set of bits the number of bits affected depends on the data rate and duration of noise for example, if we are sending data at 1kbps, a noise of 1/100s can affect 10bits, of we are sending data at 1mbps, the same noise can affect 10,000bits

## 2.36   Correction verses Detection

**Correction**: The correction of errors is more difficult than detection. In error detection, we are looking only to see if any error has occurred. The answer is simple Yes or No. we are not even interested in the number of errors. A single-bit error is the same for us as a burst error.

In error correction, we need to know the exact number of bits that are corrupted and more importantly, their location in the message the number of errors and the size of the message are important factors.

If we need to correct one single error in a 8-bit data unit, we need to consider eight possible error locations; if we need to correct two errors in a data unit of the same size, we need to consider 28 possibilities.so imagine the receiver's difficulty in finding 10 errors in data unit of 1000 bits.

## 2.37   Forwarding ERROR Correction versus Retransmission

Correction is the process in which the receiver tries to guess the message by using redundant bits. This is possible, as we see later, if there are two main methods of error correction. Forward error number of error is small correction by retransmission is a technique in which the receiver detects the occurrence of an error and asks the sender to repeated until a message arrives that the receiver believes is error-free (usually, not all errors can be detected).

## 2.38   Redundancy

The central concept in detecting or correcting errors is redundancy to be able to detect or correct errors, we need to send some extra bits with our data. These redundant bits are added by sender and remove by the receiver in the destination. Their presence allows the receiver to detect and correct errors we must and we need to send extra (redundant) bits with our data.

## Coding

Redundancy is achieved through various coding sachems. The sender adds redundant bits through a process that creates a relationship between the redundant bits and the actual data bits. The receiver checks the relationship between the two sets of bits to detect or correct the errors. The ratio of redundant bits to the data bits and the robustness of the process are important factors in any coding scheme the general idea of coding to us and we can divide coding schemes into two broad categories: block coding and convolution coding. In here I will concentrate on block coding as considering my bachelor degree: but convolution coding is more complex and beyond the scope of my thesis.

## 2.39   Modular Arithmetic

Before we finish this section, let me briefly discuss a concept basic to computer science in general and to error detection and correction in particular: modular arithmetic our intent here is not to delve deeply into mathematics of these topics; we present just enough information to provide a background to materials discussed in this chapter.

In modular arithmetic's, we use only a limited range of integers we define an upper limit; called a modulus N. we then use only the integers 0 to N-1, inclusive this is module 0-N arithmetic. For example, if the modulus is 12, we use only the integers 0 to 11, inclusive. An example of modulo arithmetic is our clock system. It is based on modulo – 12 arithmetic substituting the number 12 for 0. In a modulo 0-N system, if a number is greater than N, it is divided by N and the remainder is the result. If it is negative, as many (Ns) as needed are added to make it positive. Consider our clock system again, if we start job at 11 am and the job takes 5h, we can say that the job is to be finished at 16:00 if we are in the military, or we can say that it will be finished at 4 P.M (the remainder of 16 h is 4) So in modulo –N arithmetic, we use only the integers in the range 0 to N-1, inclusive.

Addition and subtraction is no carry when you add two digits in column. There is no carry when you subtract one digit from another in a column.

## 2.40   Block coding

In block coding, divide our message in to blocks, each of Kbits, called data words. We add R redundant bits to each block to make the length n=k+r. The resulting n-bit blocks

are called code words. How the extra r bits is chosen or calculated is something we will discuss late. For moment it is important to know that we have a set of data words, each of size k, and a set of code words, each of size N. with k bits, we can create a combination of 2power k data words; with n bits, we can create a combination of 2 power n code words since n>k, The number of possible code words is larger than the number of possible data words. The block coding process is one –to-one; the same data word is always encoded as the same code word.

## 2.41   Error detection

How errors can be detected by using block code?

If the following two conditions are met, the receiver can detect a change in the original code word.

1.      The receiver has (or can find) a list of valid code words.

2.      the original code word has changed to an invalid one

Show error detection in block coding. And it is the process of error detection in block coding the sender creates code words out of data words by using a generator that applies the rules and procedures of encoding. Each code word sent to the receiver may change during transmissions. If the received code word is the same as one of the valid code words, the word is accepted; the corresponding data word is extracted for use. If the received code word is not valid, it is discarded. However if the code word is corrupted during transmission but the received word still matches a valid code word, the error remains undetected this type of coding can detect only single errors. Two or more errors may remain undetected.

Assume the sender encodes the data word (01) as (011) and send to the receiver consider the following cases:

1.      The receiver receives 011. It is valid code word. The receiver extracts the data word 01 from it.

2.      The code word is corrupted during transmission and 111 is received (the left most bits is corrupted). This is not a valid code word and is discarded.

3.      The code word is corrupted during transmission, the 000 is received (the right

two bits are corrupted). This is a valid code word the receiver incorrectly extracts the data words 00. Two corrupted bits have made the error undetectable.(An error-detecting code can detect only the types of errors for which it is designed; other types of errors may remain undetected.

## 2.42   Congestion Control

In here, we discuss various techniques for controlling congestion in packet switching, frame relay, and ATM networks, and ip based internets. To give context to this discussion provides a general depiction of important congestion control techniques.

## 2.43   Backpressure

We have already made reference to backpressure as a technique for congestion control. This technique produces an effect similar to backpressure in fluids flowing down a pipe. When end of the pipe to the point is closed (or restricted), the fluid pressure backs up the pipe to the point of origin, where the flow is stopped (or slowed).

Backpressure can be exerted on the basis of links or logical connections (e.g., virtual circuits). Referring again to if node 6 becomes congested (buffers fill up), then node 6 can slow down or halt the flow of all packets from node 5 (or node 3, or nodes 5and 3). If this restriction persists, node 5 will need to slow down or halt traffic on its incoming links. This flow restriction propagates backward (against the flow of the data traffic) to sources, which are restricted in the flow of new packets in to the network.

Backpressure can be selectively applied to logical connections, so that the flow from one node to the next is only restricted or halted on same connections, generally the ones with the most traffic. In this case, the restriction propagates back along the connection to the source.

Backpressure is limited utility. It can be used in a connection –oriented network that allows hop-by-hop (from one node to the next ) flow control. X.25-based packet-switching networks typically provide this feature. However, neither frame relay nor ATM has any capability for restricting flow on a hop-by-hop basis. In the case of IP – based internets, there have traditionally been no built-in facilities for regulating the flow of the data from one router to the next along a path through the internet. Recently, some flow-based schemes have been developed; this topic is introduced in part five.

## 2.44  Choke Packet

A choke packet is a control packet generated at a congested node and transmitted back to a source node to restrict traffic flow. An example of a choke packet is the ICPM (Internet Control Message Protocol) Source Quench packet. Either a router or a destination end system may send this message to a source end system, requesting that reduce the rate at which it is sending traffic to the internet destination. On receipt of a source quench message, the source host should cut back the rate at which it is sending traffic to the specified destination until it no longer receives source quench messages. The source quench message can be used by a router or host that must discard IP datagrams because of a full buffer. In that case, the router or host will issue a source quench message for every datagram that it discards. In addition, a system may anticipate congestion and issue source quench messages when its buffers approach capacity. In that case, the datagram referred to in the source quench message may well be delivered. Thus, receipt of a quench message does not imply delivery or no delivery of the corresponding datagram.

The choke package is a relatively crude technique for controlling congestion

Some others techniques of congestion controlling

## 2.45  Implicit congestion signaling

When network congestion occurs, two things may be happen (1) The transmission delay for an individual packet from source to destination increases, so that it is noticeably longer than the fixed propagation delay, and (2) packets are discarded. If a source is able to detect increase d delays and packet discards, then it has implicit evidence of network congestion. If all sources detect congestion and in response reduce flow on the basis of congestion, then the network congestion will be relieved. Thus, congestion control on the basis of implicit signaling is the responsibility of end systems and does not require action on the part of network nodes.

Implicit signaling is an effective congestion control technique in connectionless, or datagram, configurations such as datagram packet-switching networks and IP based internets. In such cases, there are no logical connections through the internet on which flow can be regulated. However, between the two end systems logical connections can be established at the TCP level. TCP includes mechanisms for acknowledging receipt of

TCP segments and for regulating the flow of the data between source and destination on a TCP connection. TCP congestion control techniques base d on the ability to detect increased delay and segment loss are discussed in chapter 20Implicit signaling can also be used in connection-oriented networks. For example, in frame relay networks, the LAPF control protocol, which is end to end, includes facilities similar to those of TCP for flow and error control. LAPF control is capable of detecting lost frames and adjusting the flow of data accordingly.

## 2.46 Explicit congestion signaling

It is desirable to use as much of the available capacity in a network as possible but still react to congestion in a controlled and fair manner. This is the purpose of explicit congestion avoidance techniques. In general terms for explicit congestion avoidance, the network alerts end systems to growing congestion within the network and the end systems take steps to reduce the offered load to the network.

Typically, explicit congestion control techniques operate over connection-oriented networks and control the flow of packets over individual connections. Explicit congestion signaling approaches can work in one of two directions.

## 2.47 Backward

Backward notifies the source that congestion avoidance procedures should be initiated where applicable for traffic in the opposite direction of the received notification. It indicates that the packets that the user transmits on this logical connection may encounter congested resources. Backward information is transmitted either by alerting bits in header of a data packet headed for the source to be controlled or by transmitting separate control packets to the source.

## 2.49 Forward

Forward notifies the user that congestion avoidance procedures should be initiated where applicable for traffic in the same direction as the received notification. It indicates that this packet, on this logical connection, has encountered congested resources. Again, this information may be transmitted either as altered bits in data packets or in separate control packets. In some schemes, when a forward signal is received by an end system, it echoes the back along the logical connection to the source. In other schemes, the end system is expected to exercise flow control upon the source end system at a higher layer

(e.g. TCP).

**Explicit congestion Categories**

We can divide explicit congestion signaling approaches into three general categories:

**Binary**

A bit is set in a data packet as it is forwarded by the congested node. When a source receives a binary indication of congestion on logical, it may reduce its traffic flow. Connection

**Credit based**

These schemes are based on providing an explicit credit to a source over a logical connection. The credit indicates how many octets or how many packets the source may transmit. When the credit is exhausted, the source must await additional credit before sending additional data. Credit-based schemes are common for end-to-end flow control, in which a destination system uses credit to prevent the source from overflowing the destination buffers, but credit-based schemes have also been considered for congestion control.

**Rate based**

These schemes are based on providing an explicit data rate limit to the source over a logical connection. The source may transmit data at a rate up to the set limit. To control congestion, any node along the path of the connection can reduce the data rate limit in a control message to the source.

## 2.50  Traffic Management

There are number issues related to the congestion control that might include under the general category of traffic management. In its simplest form, congestion control is concerned with efficient use of a network at high load. The various mechanisms discussed in the previous section can be applied as the situation arise, without regard to the particular source or destination affected. When a node is saturated and must discarded packets, it can apply some simple rule, such as discarded the most recent arrival. However, other consideration can be used to refine the application of congestion control techniques and discard policy. We briefly introduce several of those areas here.

**Fairness**

As congestion develops flow of the packets between sources and destination will experience increased delays and with high congestion. Packet losses in the absence of other requirement, we would like to assure that the various flows suffer from congestion equally. Simply to discard on a last-in-first-discarded basis may not fair. As an example of a technique that might promote fairness, a node can maintain a separate queue for each logical connection or for each source-destination pair. If all of the queue buffers are equal length, then the queues with the highest traffic load will suffer discards more often, allowing lower-traffic connections a fair share of the capacity.

## 2.51 Quality of Service

We might wish to treat different traffic flows differently. For example, some applications, such as voice and video, are delay sensitive but loss insensitive. Others, such as file transfer and electronic mail, are delay insensitive but loss sensitive. Still others, such as interactive graphics or interactive computing applications, are delay sensitive and loss sensitive. Also different traffic flows have different priorities; for example, network management traffic, particularly during times of congestion or failure, is more important than application traffic.

It is particularly important during periods of congestion that traffic flows with different requirements be treated differently and provided a different quality of service (Qos). For example, a node might transmit higher-priority packets ahead of lower-priority packets in the same queue. Or a node might maintain different queues for different Qos levels and give preferential treatment to the higher levels.

## 2.52 Reservations

One way to avoid congestion and also to provide assured service to applications is to use a reservation scheme. Such is an integral part of ATM networks. When a logical connection is established, the network and the user enter into a traffic contract, which specifies a data rate and other characteristics of the traffic flow. The network agrees to give a defined Qos so long as the traffic flow is within contract parameters; excess traffic is either discarded or handled on a best-effort basis, subject to discard. If the current outstanding reservations are such that the network resources are inadequate to meet the new reservation, then the new reservation is denied. A similar type of scheme

has now been developed for IP-based **internets.**

One aspect of reservation scheme is traffic policing. A node in the network, typically the node to which the end system attaches, monitors the traffic flow and compares it to the traffic contract. Excess traffic is either discarded or marked to indicate that it is liable to discard or delay.

# CHAPTER III
# REQUIREMENT GATHERING AND ANALYSIS

## 3.1 Auditing

Auditing is the important issue in tracking, traffic management, congestion control, error detection and its correction for better transmission is the important task to make the best network system. Some other technique may also include for a better network like using some protocols for available and integral system.

All the network systems for their high availability uses some redundant devices to be survive for services.

## 3.2 Alternate devices for redundancy

Redundancy waits for a specific event or any kinds of deactivation of some systems to occur, and then the redundant device will operate instead and take the operation. Deactivation events in in the system can be divided into many types and they are by devices failure and can be by power supplies. Redundancy of power supplies can be the alternate generators and UPS for redundancy of the network devices can be the alternative routers bridges hubs and switches and those may cost a little more but high integrity and high availability for customer services.

## 3.3 Load balancers

Load balancing is the most straightforward method of scaling out an application server infrastructure. As application demand increases, new servers can be easily added to the resource pool and the load balancer will immediately begin sending traffic to the new server. Core load balancing capabilities include:

Layer 4 (L4) load balancing - the ability to direct traffic based on data from network and transport layer protocols, such as IP address and TCP port

Layer 7 (L7) load balancing and content switching – the ability to make routing decisions based on application layer data and attributes, such as HTTP header, uniform resource identifier, SSL session ID and HTML form data

Global server load balancing (GSLB) - extends the core L4 and L7 capabilities so that they are applicable across geographically distributed server farms

How do load balancers work?

When one application server becomes unavailable, the load balancer directs all new application requests to other available servers in the pool.

To handle more advanced application delivery requirements, an application delivery controller (ADC) is used to improve the performance, security and resiliency of applications delivered to the web. An ADC is not only a load balancer, but a platform for delivering networks, applications and mobile services in the fastest, safest and most consistent manner, regardless of where, when and how they are accessed.

## 3.4    Techniques are used for detection and correction of errors

Error detection and correction or error control is technique that enables reliable delivery of digital data over unreliable communication channels. Many communication channels are subject to channel noise, and thus errors may be introduced during transmission from the source to a receiver. Error detection techniques allow detecting such errors, while error correction enables reconstruction of the original data in many cases some techniques are mostly used in detection and correction of data are redundancy in bits and retransmission of the data.

## 3.5    Techniques for congestion control

Congestion control refers to the techniques used to control or prevent congestion. Congestion control techniques can be broadly classified into two categories:

**Open Loop Congestion Control**: in here the techniques are as follow.

## 3.6    Retransmission Policy:

It is the policy in which retransmission of the packets are taken care. If the sender feels that a sent packet is lost or corrupted, the packet needs to be retransmitted. This transmission may increase the congestion in the network.

## 3.7    Window Policy:

The type of window at the sender side may also affect the congestion. Several packets in the Go-back-n window are resent, although some packets may be received successfully at the receiver side. This duplication may increase the congestion in the network and

making it worse.

## 3.8 Discarding Policy:

A good discarding policy adopted by the routers is that the routers may prevent congestion and at the same time partially discards the corrupted or less sensitive package and also able to maintain the quality of a message.

## 3.9 Admission Policy:

In admission policy a mechanism should be used to prevent congestion. Switches in a flow should first check the resource requirement of a network flow before transmitting it further. If there is a chance of a congestion or there is a congestion in the network, router should deny establishing a virtual network connection to prevent further congestion.

## 3.10 Closed Loop Congestion Control:

Backpressure, Choke Packet Technique, Implicit Signaling, Explicit Signaling.

## 3.11 Traffic Management Techniques for Packet Networks

Regardless of the technology used, certain fundamental traffic management techniques can be applied to optimize utilization and handle congestion in packet-switched networks and some technique are as follow:

## 3.12 Classification:

It is the process of determining for each packet entering a network or node which of several defined classes it belongs to. Classification may be based on existing QoS markings, or on more general parameters of the packet, such as source address, destination address, or application being carried.

## 3.13 Marking:

It is the process of placing a QoS value into the packet header, so that downstream processes n this node or other nodes can identify the QoS of the packet correctly and handle its queuing/scheduling and discarding appropriately

## 3.14 Policing:

It is normally carried out at the head-end of a connection or flow to ensure that the contracted traffic rate is not being exceeded. Nodes may use shaping to smooth excessive bursts of traffic, and may discard or remark traffic if necessary according to the policies of the contract which applies

## 3.15 Congestion Management:

It deals with situations where more traffic than can be carried by a node is present to the output ports.

## 3.16 Queuing and Scheduling:

It allows different handling of packets based upon their class of service as defined by header markings or other classification within a node. Simple priority scheduling services the most important queue until it is empty, then services the next priority queue in the same way. Most nodes operate some more sophisticated weighted queuing algorithm which ensures some bandwidth for all priority classes

# CHAPTER IV
# METHODOLOGY

## 4.1 Area of study

This study was started in Kabul Afghanistan, the city I live in and there are most of the telecommunication companies ISP companies some other business companies all those companies need for better network and optimized network systems for better advantages to gain from customers and provide the best service for them.

Most of the public and private sectors are also attracted to expand their network and they are trying to reduce the latency in the network and keep their network secure to provide better services for the public.

## 4.2 Sampling

The study includes some ISP companies, IT departments, and some parts of telecommunication camponies, Local companies IT staff and some other organizations staffs.

The researcher used purposive stratified and random sampling methods to select the sample for the study.

The study used simple random sampling design. The basic idea is that each of companies if it is telecom or ISP and public, private sectors had an equal chance of being included in the sample.

Purposive sampling was use because the sample wanted was already known.

## 4.3 Research Design

Research design is the methodology or it is the procedure as a plan of action through which a researcher organizes his /her work from data collection, organization to analysis.

This study used the case study design. The goal was to have the detailed information about the optimization of the network and also to prevent the disappearance of the system in causes of network down timing.

## 4.4    Research Strategy

The qualitative research strategy was used to this study. Qualitative is scientific method of observation to gather non numerical data. This type as refers to the meanings, concepts, definitions, and other some characteristics. The nature of the study problem requires much of listening and observation hence qualitative research was appropriate to this study.

## 4.5    Data Collection Methods

Both primary and secondary collected for this study. The techniques for data collection method, which are applicable to the qualitative research process, were therefore use. Secondary data techniques used include:

## 4.6    Secondary

I have collected data from web pages and secondary method like research papers, study from books video learning's focus group and extra interviews methods in-depth interviews IT director of some networked companies like ISP company one of them.

## 4.7    Focus Group Discussion

The focus group discussion was conduct because this method of data collection enables to see how people respond to each other's view and build up views out of the interactions that take place within the group. The method created for getting diversity in perspectives.

With the focus group, the network optimization and prevention from latency and loss of the network was discussed in focus group.

## 4.8    Observation

Direct observation was also used to this study as a source of obtaining information. The activities of these two sectors public and private also observed during the study.

This enabled getting personal experience and the depth of knowledge on these two sector technologies and its impact to sustainable development, Impact on private organization, Impact on the public sectors. The method was also useful as the field of study was still largely unexplored.

## 4.9 Data analysis methods:

The study used qualitative methods to analyze the collected data.

Tables and figures present the results. The researcher combined the information from all sources of information that is interview, group discussions, observation, and secondary sources means that the data is collected from web blogs and other researcher's secondary sources.

## 4.10 Limitation

This study is an attempt to analyze the impact some optimized system in the network and to avoid the latency from the system.

# CHAPTER V
# IMPLEMENTATION OF PROJECT

## 5.1    Fault tolerant network

A design of fault tolerant enables a system to continue its intended operation, possibly at a reduced level rather than failing completely when some part of the system fails. The term is most commonly used to describe computer systems designed to continue more or less fully operational with perhaps, a reduction in through put or an increase in response time in the event of some partial failure. That is the system as a whole is not stooped due to problems either in the hardware or the software. An example in another field is a motor vehicle designed so it will continue to be drivable if one of the tires is punctured, or a structure that is able to retain its integrity in the presence of damage due to causes such as fatigue, corrosion, manufacturing flaws, or impact. In the figure we can see the fault tolerant network with the redundant devices for high availability.
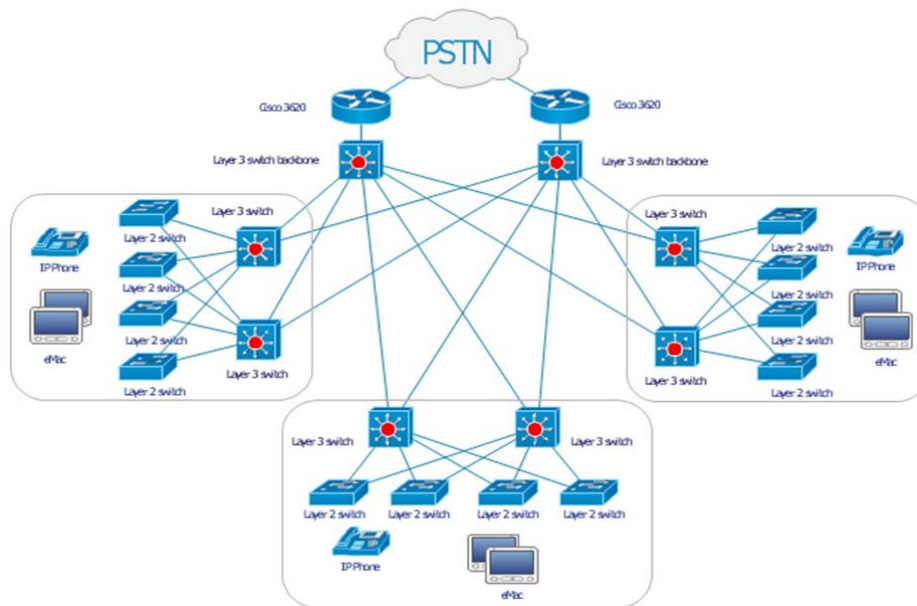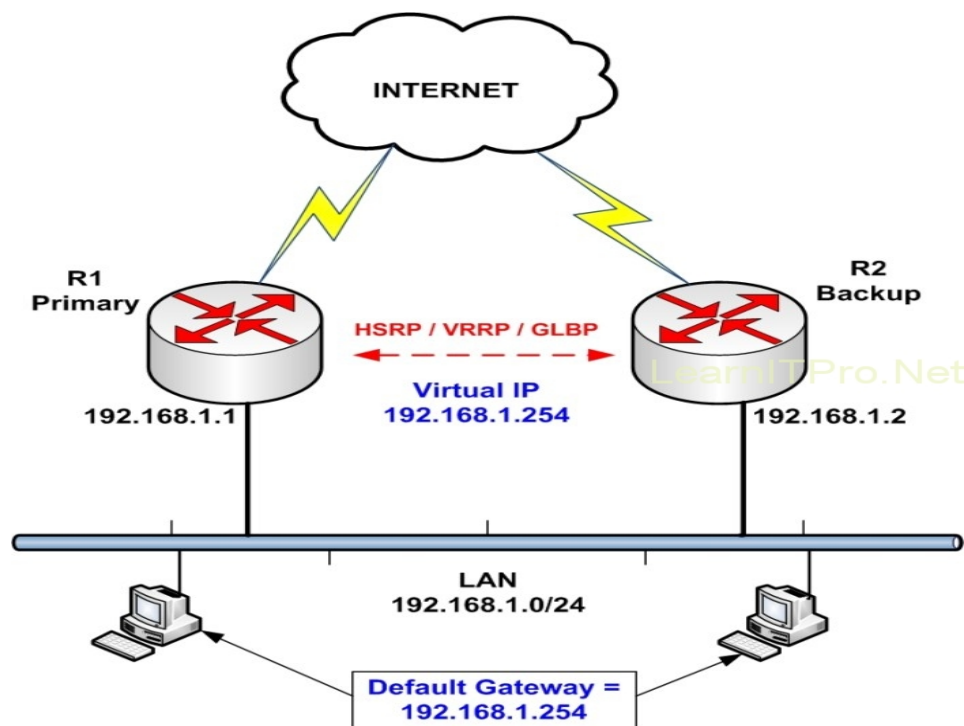


Figure  10 -FTN

## 5.2    Network Redundancy protocols

Network redundancy is the process through which additional or the alternate instances of the network devices, equipment, protocols and communication mediums are installed within the network infrastructure. Redundancy protocols are the standards for that to provide seamless failover against failure of any network component. These redundancies are invisible to the application. And these protocols are installed in tow nodes of ports and are attached to two separated networks of similar topology. This method is for the ensuring network for a high availability in case of avoiding the network devices or path failure and unavailability. Network redundancy is primarily implemented in enterprise networks infrastructure to provide a redundant source of network communications. It will serves as a backup mechanism for quickly swapping network. In this figure the HSRP, VRRP, GLBP will configure.



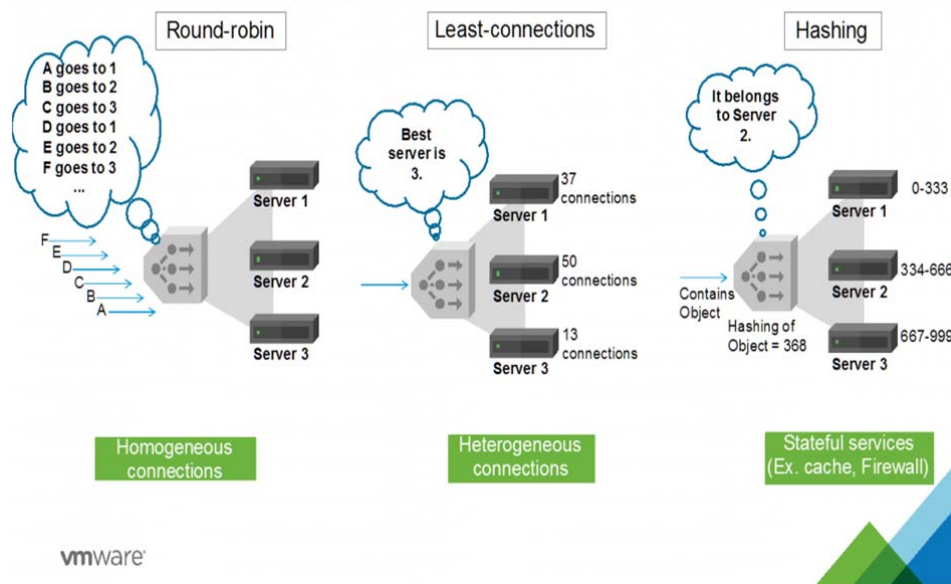## 5.3    Load Balancing and its Algorithms

As we know load balancing is the process of dividing the load to multiple servers so a load balancer acts as the "traffic cop" sitting in front of your servers and routing client requests across all servers capable of fulfilling those requests in a manner that maximizes speed and capacity utilization and ensures that no one server is overworked,

which could degrade performance. If a single server goes down, the load balancer redirects traffic to the remaining online servers. When a new server is added to the server group, the load balancer automatically starts to send request to it.

Load balancing possess of three algorithms

- Round robin the first one in the figure

- Least-connections the second one in the figures

- Hashing is the third one in the figure
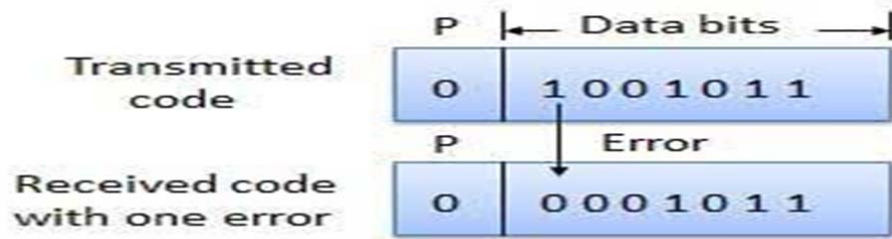


**Basic Load Balancing Algorithms**

## 5.4    Error Detection and correction

- Error detection refers to a class of techniques for detecting garbled message
Adding some extra bits to detect occurrence of error
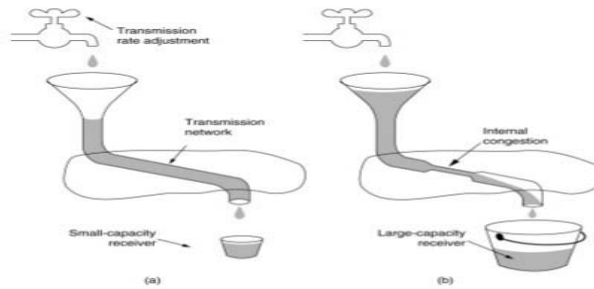Not enough to detect the position of errors.

- Error correction adding enough redundant bits to reduce what the correct bits are
Error correction is too expensive and hard.

|  | P | ← Data bits → |
|---|---|---|
| Transmitted code | 0 | 1 0 0 1 0 1 1 |
|  | P | Error |
| Received code with one error | 0 | 0 0 0 1 0 1 1 |

## 5.5    Congestion control

Congestion control window controls the volume of packets in the flight

- Slow starts/ congestion avoidance

- Retransmit time out

- Fat retransmit and fast recovery



(a) A fast network feeding a low capacity receiver.
(b) A slow network feeding a high-capacity receiver.

# CONCLUSION

It can be a big problem that when our system goes down and there won't be a backup or redundant devices for taken the operation so we need an optimized network with redundant devices and machines to don't faces any problems in a network system the optimized network can include Fault tolerance, Redundancy, Load balancing, Error detection and its correction, Congestion control, traffic management so with these techniques we can optimize our network system in best way to have the best operation in our progress and decrease the response time and reduce the latency in the network.

Redundant devices can be the Network devices power supply and generators so if we have the redundancy in the network high availability will come too for example to have two generators one fails the redundant one will take operation instead .

## 6.1 Network Redundancy protocols

Are the protocols used to make the high availability in the systems as we have the many Network redundancy protocols like HSRP, VRRP, GLBP, CARP, STP and so on each of those protocols possess of its own properties and its function all them work to make one device active another device standby mode if active fail the standby mode machine will take the operation.

## 6.2 Load balancing

It is the another technique of making the network available and increase the response time and reduces latency for the system and user in the networks the load balancing is the process of the distributing the load of request data or the on many different machines and do the functions so fast then the response time will be less and the latency will not occur if occurred it will so less because the data is divided into many machine if one of them don't worked the request will direct the machine or the server which finished its task and it will pick up another task too.

# References

[1] A study on network optimization problems and its solutions by Der-San Chen, Robert G, Batson, and Yu Dang Authors.

[2] A fault tolerant Engineered Network A study on fault tolerance Authors Vincent Liu, Daniel Halperin, Arvind Krishnamurthy, and Thomas Anderson, University of Washington.

[3] A study on network latency and how to decrease it by Boris Rogier from.

[4] Algorithm for network Redundancy study on its Advantages by Yong Xu from Southwest university of Science and Technology, Mian Yang, China.

[5] Data and computer communications by William Stallings Study on optimizing a network flow with error control congestion control traffic management.

[6] Performance Evaluation of FHRP (HSRP,VRRP,GLBP) By Zia Ur Rahman, Safyan Mukhtar, sajjad khan, raees khan, Reena Rashid, waqas khan from Department of Computer science bacha khan University, Chaharsada Pakistan.

[7] Reliability of computer systems and Networks with fault tolerance by Martin L.shooman.

[8] Server Load balancing By Tony Bourke from.

[9] The importance of Network Redundancy by Robert Shimonski from.

[10] Princewill Aigbe and Jackson Akpojaro, "Analysis of security issues in Electronic Payment System, international journal of computer application (0975-8887), 10, December 2014.

[11] S.Sumanjeet, "Emergence of Payment Systems in the age of Electronic Commerce: The State of Art, Asia Pacific Journal of Finance and Banking Research.

[12] Goyal, "Mobile Banking in India: Practices, challenges and security issues, international journal of Advanced Trends in Computer Science and Engineering, 2012.

[13] 24 April 2015. [Online]. Available: http://www.theengineer.co.uk/more-sectors/electronics-chequebookcoul-.

[14] Niranjanamurthy, "The study of E-Commerce Security Issues and Solution." International journal of advanced research computer and communication engineering, vol. 2, no. 7, 2013.

[15] Uddin, "E-Wallet System for Bangladesh an Electronic Payment System".

[16] http://www.ofnisystems.com/services/validation/computer-systems/.

[17] Test. [Online]. Available: FDA 21 CER 11.10(a) and EMA Annex 11, section.

[18] Information systems project management: methods, tools and techniques / John McManus and Trevor Wood-Harper, Harlow [etc.] : Prentice Hall, c2003

[19] Subversion version control: using the Subversion version control system in development projects / William Nagel, Upper Saddle River (N.J.): Prentice Hall/PTR, c2005

[20] Systems Analysis and Design Shelly Cashman Adamski Boston 1991

[21] Software Engineering Roger S.Pressman UK, c2000, 5th ed.

[22] Ubuntu - http://packages.ubuntu.com/intrepid/libapache2-mod-wsgi

[23] Open SSH - http://sial.org/howto/openssh/publickey-auth/