# Reliability and Performance Analysis of New Fault Tolerant Irregular Network

RITA MAHAJAN,
Department Of Electronics and Electrical Communication,
Punjab Engineering College, Chandigarh,
INDIA.
rita_mahajan@rediffmail.com

RENU VIG,
Department of Electronics,
UIET, Panjab University, Chandigarh,
INDIA.
renuvig@hotmail.com

*Abstract:*    A lot of work and analysis has been done on regular fault-tolerant MINs but irregular MINS have the inherent concept of favorite memory module offering lower latency. The multipath nature of these networks provides a non-blocking operation. This work introduces a new network, named as MFT-2 network, which contributes fault-tolerance by providing dynamic full-access capability in the presence of faults. MFT-2 network is improvement on irregular Four Tree network (FT). In case of any fault the system  is reconfigured which  operates in a degraded mode owing to the increased latency. But this approach minimizes the overheads of providing fault-tolerance, both in terms of cost and performance, during normal operation of the system. Irregular nature of these networks has reduced latency, improved dynamic fault-tolerant routing, high reliability and performance. The probabilistic approach is used to analyze the MINs based on independent request assumptions.

*Keywords*: - Multistage Interconnection Networks (MIN), Irregular Network, Parallel Processing, DOT Network, FT (Four-tree network), Reliability.

## 1    Introduction

In recent years, parallel and distributed computing has become very popular as it provides a means to overcome the limitations imposed by sequential computer [4]. A common situation in parallel system is one in which set of processor have to access a set of memory modules or processor [5]. Vital components of these systems are the Interconnection Network that enables the processors to communicate among them [1][3]. The performance of multi-processor system depends primarily on design of its Interconnection Network [2][7]. Multistage Interconnection Network (MIN) achieves the full access property with far fewer connections and hence become popular in short duration of time. A MIN consists of a stack of switching element (SE) connected with permuters. The MIN are designed for N input and N output using m x m SE's. There are two types of MIN: regular and irregular network[8]. A regular MIN contains same number of switching element at each stage of network whereas the number of switching element in different stages may differ in an irregular network [19]. A new irregular fault tolerant MIN called MFT-2 is proposed and analyzed. MFT-2 network achieve the goals of a fault-tolerant network i.e. well performance even

in presence of fault, high Permutation possibility, less average path length, low cost and simple control scheme. The paper is organized as follows: Section 2 provides a brief introduction to previously proposed networks like DOT, MDOT and FT networks. Section 3 describes the construction of MFT-2 network. Section 4 discusses the path length, routing tag algorithm and routing procedure. Section 5 discusses the fault-tolerance of MFT-2. Reliability of MFT-2 has been evaluated in terms of MTTF in section 6. Various performance parameters have been compared in section 7. In section 8 cost effectiveness has been discussed. Finally conclusion is discussed in section 9.
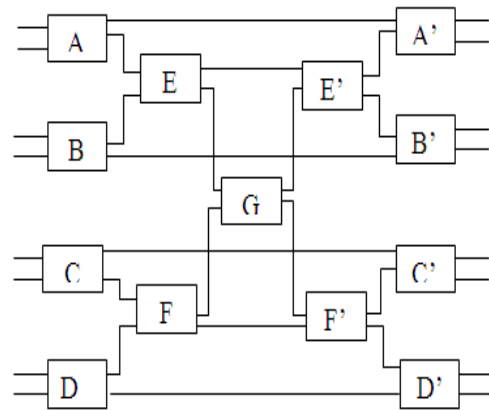


Fig. 1: Double Tree network

## 2    Previously Proposed Networks

A DOT network (Shen 1982) is an irregular type of multistage interconnection network. It consists of a right and a left half, each resembling a binary tree [20][15]. The two halves are mirror images of each other. A DOT network of size $2^n$ x $2^n$ has $2^n$ source and $2^n$ destination terminals and 2n-1 stages. Further, it has $2^{n+1}$ –3 switches. An $i^{th}$ and $(2n-i)^{th}$ stage has $2^{n-i}$ switches of size 2x2 for i=1,2,3…n. Fig. 1 shows a 8 x8 DOT network[11]. The flip control (i.e. individual stage control) used in the DOT network affects the performance and reliability of the entire system[21]. The connections between the switches of a DOT network has been slightly changed thus the resulting network as shown in Fig. 2 is termed as modified Double Tree (MDOT) network[13]. This has the advantage of distributed control (i.e. individual switch control). The MDOT network is used in the construction of Four Tree network. A FT network [10] of size $2^n$ x $2^n$ is constructed with the help of two identical groups, each consisting of a MDOT network of size $2^{n-1}$ x $2^{n-1}$, which are arranged one above the other. It consists of (2m-1) stages and $(2^{m+2} -6)$ switches where   m = $\log_2$ N/2 and N=$2^n$.
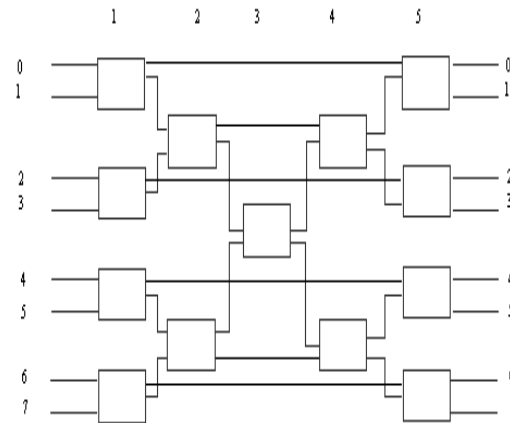


Fig. 2: Modified Double Tree network

A FT network has $2^{n-1}$ switches of size 2x2 and rest of size 3x3. There are $2^n$ multiplexer of size 2x1 and $2^n$ demultiplexer of size 1x2. Both stage i and stage (2m-i) have exactly $2^{n-i}$ switches where i = 1, 2…2m-1. The switches in a stage having the same number in both groups form a loop. Such loops of switches are formed in all stages except the last one. The cost complexity of FT is (9.75 $2^{n+1}$ -54). The FT network is shown in Fig. 3.
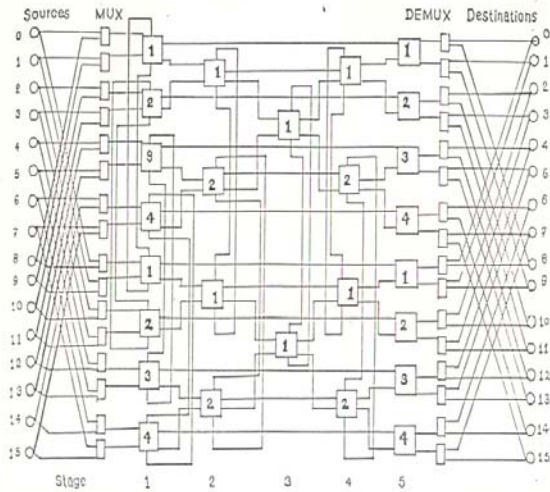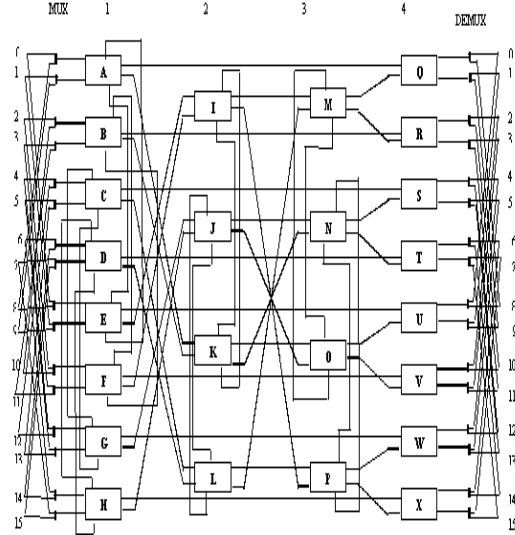
**Fig. 3: Four-Tree Network**



**Fig. 4: MFT- 2 (Proposed network)**

## 3    Construction procedure of MFT-2 network

MFT-2 network of size $2^n$ x $2^n$ is constructed with the help of FT network as shown in Fig 4. MFT-2 is constructed by removing the middle stage from FT network, which is a source of blocking. The new MFT-2 network of Size $2^n$ x $2^n$ consist of 2m-2 stages where m= $\log_2$ N/2 and N= $2^n$. It has $2^{m+2}$ – 8 switches out of which $2^{n-1}$ are of 2 x2 where as rest are of 3x3. There are $2^n$ multiplexer of size 2 x 1 and $2^n$ demultiplexer of size 1 x 2. Both stage i and stage (2m-i-1) has exactly $2^{n-i}$ switches where i = 1, 2…2m-2.

## 4    Routing Scheme

The reliability and performance of a network depends upon how effectively the alternate paths available and utilized. The routing procedure starts with identification of path length and creating routing tag. The following section describes the routing scheme of MFT -2 Network.

### 4.1    Path length algorithm

For a given source- destination pair, there are multiple path of different path length in MFT-2 network. The algorithm for the allocation of path length gives the information about the different possible paths between a source and destination pair. The possible path length between source and destination vary from 2 to 2m-2 for a given $2^n$ x $2^n$ network, depending upon the source and destination terminal.

The source S and destination D is represented in binary code as

$$S = S_{n-1} S_{n-2}……….S_o$$

$$D = D_{n-1} D_{n-2}………..D_o$$

The Path length algorithm is as follow
If [$(S_{n-2} \oplus D_{n-2}) + (S_{n-3} \oplus D_{n-3}) +……………+ (S_1 \oplus D_1)$] = 0
then

Minimum path length is 2 and all the paths of different length are possible i.e. 2, 4…2m-2
Else if [$(S_{n-2} \oplus D_{n-2}) + (S_{n-3} \oplus D_{n-3}) +……………+ (S_2 \oplus D_2)$]= 0
then
All path of length equal to or greater than 4 are possible.
.
.

.
.

if $[(S_{n-2} \oplus D_{n-2}) + (S_{n-3} \oplus D_{n-3}) + \ldots \ldots \ldots \ldots + (S_j \oplus D_j) = 0$

(Where $\oplus$ represent X-OR operation, + represent OR operation)
Then all the paths of length equal to or greater than 2j are possible.
Else
Path of length 2m-2 are possible.

**Example:**
**Case 1:** Let S =0000 and D= 0001 of $2^4$ x $2^4$ MFT-2 network as $[(0 \oplus 0) + (0 \oplus 0)]$ is zero the path length is 2 and greater than 2 i.e. 2 and 4.

**Case 2**: Let S= 0000 and D= 0010 of $2^4$ x $2^4$ MFT-2 network as $[(0 \oplus 0) + (0 \oplus 1)]$ is not zero hence the path length is 4.
Table 1 shows path length from Source 0000 to all possible destinations.

**Table 1**

| S | D | Path Length |
|---|---|---|
|   | 0000 |  |
|   | 0001 | 2 or 4 |
|   | 1000 |  |
|   | 1001 |  |
|   | 0010 |  |
|   | 0011 |  |
| 0000 | 1010 |  |
|   | 1011 | Longest |
|   | 0100 | Path (4) |
|   | 0101 |  |
|   | 0110 |  |
|   | 0111 |  |
|   | 1100 |  |
|   | 1101 |  |
|   | 1110 |  |
|   | 1111 |  |

## 4.2 Routing Tag Algorithm

The routing tag algorithm for MFT-2 network gives the information about the distributed routing control required to establish a path between any source and destination terminal pair for a given path length (If it exists)

If Path length is $2 \le x < 2m-2$

Then Routing Tag = $S_{n-1} . (1.1.1 \ldots 1)_{[x/2]-1} .0 . (D_{[x/2]-1} \ldots \ldots D_0) . D_{n-1}$

Else

If x = (2m-2)
Then Routing Tag =

$S_{n-1} . (1.1 \ldots 1)_{[x/2]-1} (S_{n-2} \oplus D_{n-2}) . (D_{[x/2]-1} \ldots \ldots D_0) . D_{n-1}.$

else no tag is possible.

**Example:** For a $2^4$ x $2^4$ MFT-2 Network.

**Case 1:** Let the data be routed from Source 0 to Destination 1

S=0000                    D=0001
Path length is 2 or 4.
For x=2, Routing tag = 0.0.1.0
For x=4, Routing tag =0.1.0.0.1.0

**Case 2:** Let the data be routed from Source 4 to destination 3.

S=0100                    D=0011
Path length is 4.
For x =4, Routing tag = 0.1.1.1.1.0

**Case 3**: Let the data be routed from Source 14 to destination 4.

S=1110                    D= 0100
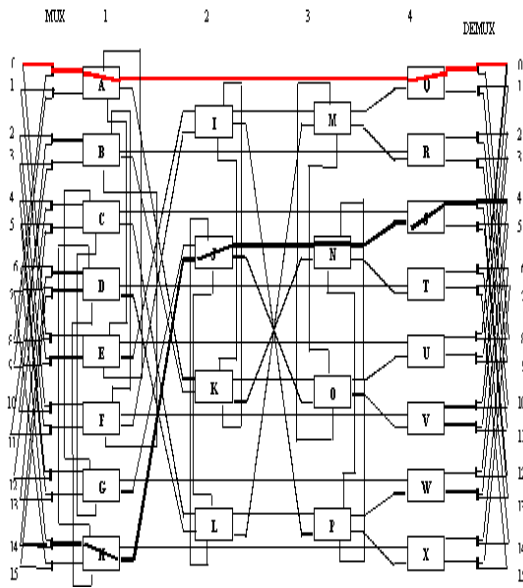Path length is 4.
For x=4, Routing tag = 1.1.0.0.0.0

**Fig.5: Routing in MFT-2 for S=14 to D=4 and S=0 to D=1**

The request is forwarded through the primary path which is of multipath length but if primary path is faulty either due to multiplexer or the switch or both then it takes the secondary. If the secondary path is also faulty then the request is dropped. The redundancy graph is shown in Fig. 6.
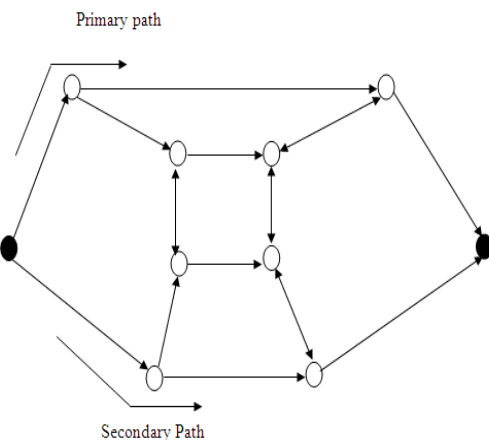


**Fig. 6: Redundancy Graph of MFT-2 Network**

# 5  Fault Tolerance

Fault tolerance is the ability of the system to continue operating in presence of fault[6]. The routing algorithm for network allows two way of routing except the final stage, which has no option to fork the request[16][17]. All other stages have a fork that helps to route through secondary path in case of faults.

**Theorem 1:** In MFT-2 network, if fault occurs at any switch of any stage except the last stage there exists at least one fault free path from any source to any destination. For Example the source 0 to destination 0 has primary path of length 2 with routing tag = (0.0.0.0) but if any fault occurs at stage 4 or a link between stage 1 and stage 4, switch A will pass the request to switch E and the request will be matured with same tag through switch E as shown in Fig. 7.

For any Source–Destination pair if the path length is four and fault occurs in stage 2 or stage 3 then the request is passed to the switch in the corresponding loop and request is matured with the same routing tag and fault is not tolerated in the stage 4.
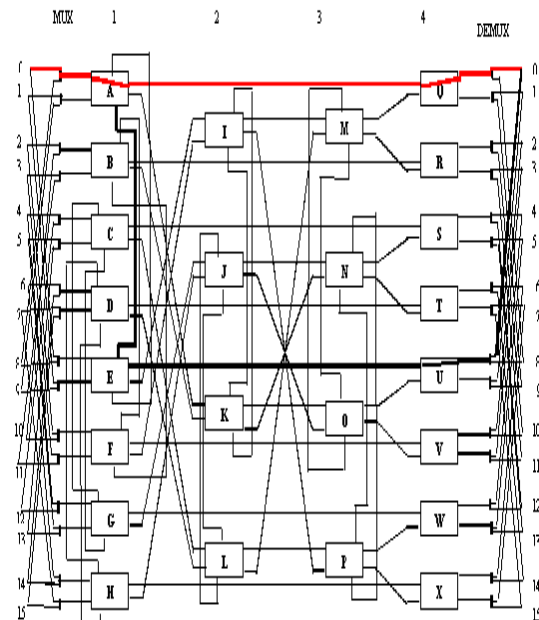


**Fig. 7: Fault tolerance of MFT-2 network**

## 6    Reliability

Dynamic MFT-2 is a multipath MIN where a redundant switching element is available at all the nodes, and the paths are used adaptively according to the faults present[9]. A series-parallel reliability model computes the bounds of reliability, accounting for degradation with time, as well as the path taken. The following assumptions are made during reliability analysis[18].

1. Switch failures occur independently in the network with the failure rate of $\lambda$, which is about $10^{-6}$ per hour.
2. Based on the gate count, failure rate of 2x2 SE is taken $\lambda_2 = \lambda$; for a 3x3 SE is taken $\lambda_3 = 2.25\lambda$ and the failure rate for a kx1 multiplexer or 1xk demultiplexer is $(k * \lambda)/4$.
3. 2x2 SEs in the last stage and their associated demultiplexers are taken as series system with a combined failure rate of $\lambda_{2d} = 2\lambda$.

**Upper Bound**

It is observed that each source is connected to two multiplexers in both the sub networks. Thus, MFT-2 is operational if one of the two multiplexers attached to a source in either sub network is operational. Reliability block diagram, shown Fig. 8(a) follows a series-parallel model. Upper Bound is computed from the following expression:



**Fig. 8(a) Upper Bound Reliability Block Diagram**

$$R_{UB\text{-}MFT\text{-}2}(t) = [1-(1-e^{-\lambda m \, t})^2]^{N/2} .$$
$$[1-(1-e^{-\lambda 3 \, t})^2]^{N/4+[N/8(2m-4)]*} .$$
$$[1-(1-e^{-\lambda 2d.t})^2]^{N/4}$$

$$MTTF = \int_0^\infty R_{UB-MFT-2}(t).dt$$

* -To be included if the path length $> 2$

**Lower Bound**

At the input side of MFT-2, routing scheme does not consider the multiplexers to be an integral part of the 2x2 switch. Hence, if both multiplexers are grouped with each switch in the input side and considered as a series system, then a conservative estimate of reliability, shown in Fig. 8(b) is obtained. Lower bound of reliability is calculated by the following expression:

$$R_{LB\text{-}MFT\text{-}2}(t) = [\ 1-(1-e^{-\lambda m \, t})^2\ ]^{N/4} .$$
$$[1-(1-e^{-\lambda 3 \, t})^2\ ]^{N/4+[N/8(2m-4)]*} .$$
$$[1-(1-e^{-\lambda 2d.t})^2]^{N/4}$$

$$MTTF = \int_0^\infty R_{LB-MFT-2}(t).dt$$

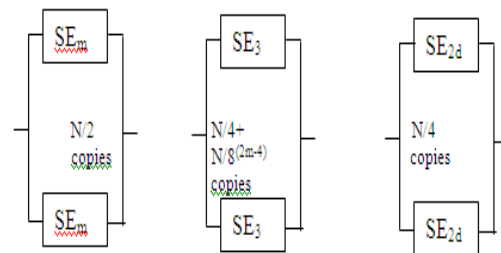* - To be included if the path length $> 2$



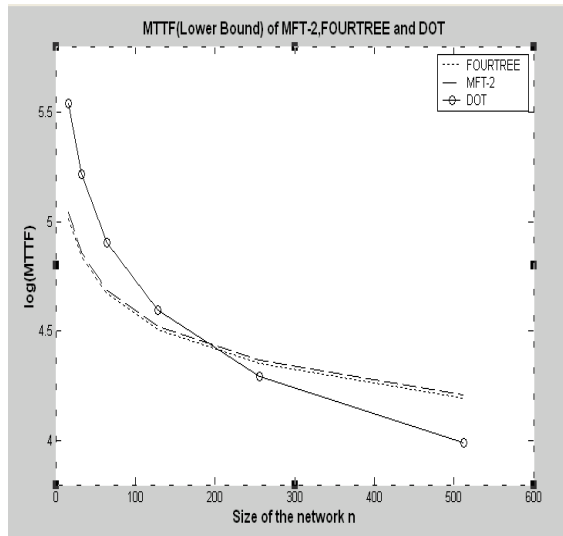**Fig. 8(b) Lower Bound Reliability Block Diagram**

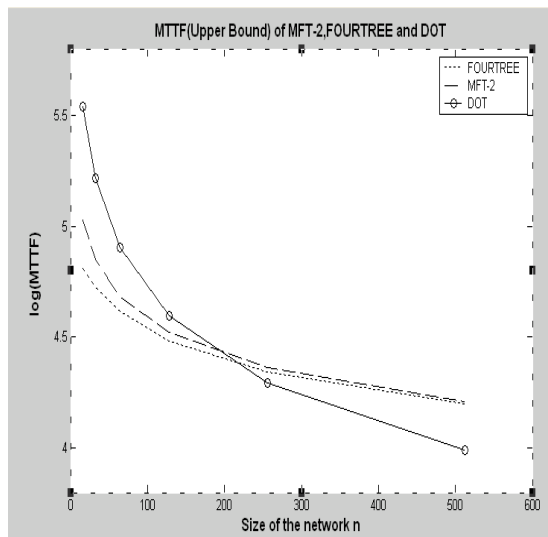**Fig. 9(a) Comparative Analysis of Lower Bound MTTF**



**Fig. 9(b) Comparative Analysis of Upper Bound MTTF**

From the Figs.9(a) and 9(b) it is seen that the lower bound MTTF and the upper bound MTTF of MFT-2 is comparable to FOUR-TREE for small network sizes and as the network size increases it shows better performance. MFT-2 is hence more reliable than the networks compared with.

As can be seen that network reliability degradation is quite less for small sizes, but increases for higher N because the path length is longer and comparatively more SEs are encountered through the route to the destination.

# 7    Performance

The analytical model has been applied to evaluate the performance of MFT-2. The performance parameters that have been evaluated are bandwidth (BW), probability of Acceptance (Pa), Throughput (TP), Processor Utilization (PU) and Processing Power (PP)[12]. Assume a MIN of size $a^n \times b^n$ constructed from a $x$ b crossbar switches and having $a^n$ Sources connected to $b^n$ Destinations. The analysis of crossbar is applied to a $x$ b crossbar switch and then extended to the complete MIN. The distinct destination digit (in base b) for setting of individual a $x$ b switches controls each stage of MIN. Since the destinations are independent and uniformly distributed, so are the destination digits [14]. For e.g. in some arbitrary stage i, a $x$ b crossbar uses digit $d_{n-i}$ of each request; this digit is not used by any other stage in the network. Moreover, no digit other than $d_{n-i}$ is used by stage i. The probabilistic approach is used to analyze the MINs based on independent request assumptions.

Given the request rate p at each of the 'a' inputs of an a $x$ b crossbar module.

The expected the rate of requests on any one of the 'b' output lines from any one input is given by:

$$\mathbf{P / b}$$

Probability of not getting the request is:

$$\mathbf{1 - p / b}$$

Probability of not getting the request from all 'a' inputs   is given by:

$$\mathbf{(1 - p / b)^a}$$

Probability of  one output getting the request from all 'a' inputs  is given by:

$$1-(1 - p / b)^a$$

The total number of requests that it passes per unit time is given by

$$b- b(1-p / b)^a$$

Dividing the above expression by the number of output lines of a x b module gives the rate of requests on any one of the b output lines.

$$1- (1 - p / b )^a$$

Thus the output rate of request q at any one of link at any stage of MIN is a function of its input rate and is given by:

$$q = 1- (1 - p / b)^a$$

Since the output rate of a stage is the input rate of the next stage, output rate of any stage can be recursively calculated starting from stage 1. And the output rate of final stage n determines the bandwidth of MIN. Let $q_i$ be the rate of request on an output link of stage i, the bandwidth for $a_n$ x $b_n$ MIN is given by:

$$BW = b^n \times p_n$$

**or**

$$BW = b^n \times p \times Pa \qquad \text{.......} \quad (1)$$

According to the definition, BW is the total number of requests matured.

$$q_i = 1 - ( 1- q_{i-1} / b )^a$$
$$q_0 = q;$$

The probability of acceptance of a request is given by:

$$P_a = b^n q_n / a^n q \qquad \text{........} \quad (2)$$

Throughput (TP):

$$TP = BW / a^n T \qquad \text{.....} \quad (3)$$

Where T is the average time taken for memory or source read / write operation.

Processor utilization (PU):

$$PU = BW / a^n p T \qquad \text{........} \quad (4)$$

Processing Power (PP):

$$PP = a^n \times PU \qquad \text{........} \quad (5)$$

These performance parameters can be computed for a general MIN model, having any number of inputs or outputs. In the following section calculations are done for N=16.

The value of $p_n$ is calculated recursively as under

$$p_0 = p$$
$$p_1 = p_0 / 2$$
$$p_2 = 1 - \prod_2 [1 - p_1 / 2]$$
$$p_3 = 1 - \prod_2 [1 - p_2 / 2]$$
$$p_4 = 1 - \prod_2 [1 - p_3 / 2]$$
$$p_n = p_4$$

The following graphs show the comparative probability of acceptance, bandwidth, throughput, processor utilization and processing power respectively.
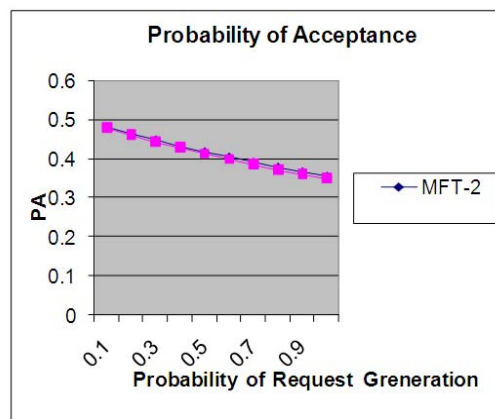


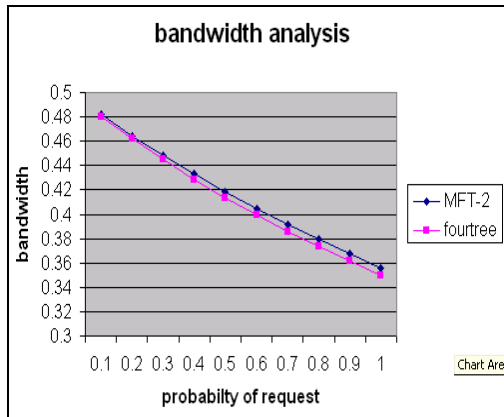**Fig. 10(a) Comparative Probability of Acceptance**
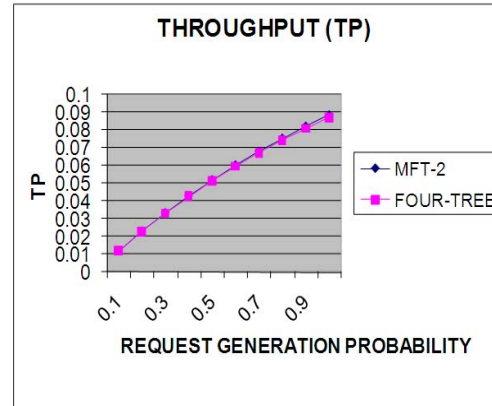
**Fig. 10(b) Comparative Bandwidth**



**Fig. 10(e) Comparative Throughput**



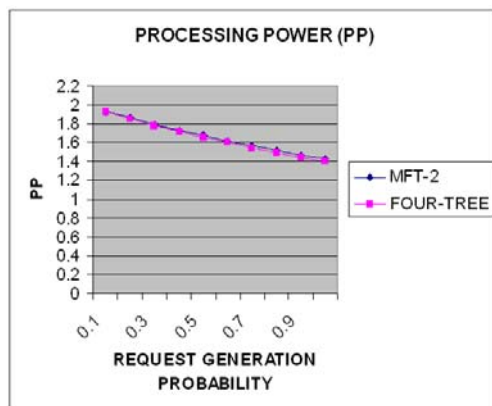**Fig. 10(c) Comparative Processor Utilisation**

The probability of generation of request (Preq-gen) affects the performance of the network. It can be seen that with the increase of this probability, Bandwidth (BW) and Throughput (TP) improves as more packets are delivered to the destination. Probability of Acceptance (Pa), Processor Utilization (PU) and Processing Power (PP) decrease because contention among switches grows with increased Preq-gen. The comparative Bandwidth, Probability of Acceptance, Throughput, Processor Utilization and Processing Power is graphically depicted in Fig. 10(a) - 10(e) respectively.



**Fig. 10(d) Comparative Processing Power**

# 8    Cost Effectiveness

Cost Complexity is quantitative term related to hardware or software or both that needs a full study of the cost to be paid in order to build and implement the network. The cost of the network can be calculated by two means: the number of connection points and number of connections or wires needed to construct the MIN.

There hardware complexity of a MIN can be calculated by two means: The number of connections points and the number of connections or wires needed to construct the MIN. Each 2x2 switch has cost of 4 units while each multiplexer of k x 1 has cost of k units. The cost complexity of a MIN is calculated in terms of cross points of all the crossbars used to built

it. The total cost complexity is the cost of each switch and links needed to build the network.

The cost of MFT-2 network depends on the number of switches in each stage. In MFT-2 network there is N 2 x 1 multiplexer and N 1x 2 de-multiplexer so net cost of multiplexer and Demultiplexer is 4N. There are N/2 switch is first stage. Each switch is of 3 x 3. There are N/4 switch is stage 2 and stage 3, each switch of 3 x 3 size then N/2 switch in fourth and last stage each stage is of 2 x 2 order. Cost complexity of MFT-2 network depends upon its multiplexer, de-multiplexer and switch at each stage of the network.

If $n = \log_r N$

(i) The input and output stage has N Multiplexer and N demultiplexer with 1 x r cross points.

(ii) The first stage has N/r SE each of (r+1) x (r+1) crosspoints.

(iii) The second and third stage has N/2r SE's each of (r + 1) x (r + 1).

(iv) The final and fourth stage has N/r SE each of (r x r).

### Net Cost Complexity

$$2Nr + N/r \, (r + 1)^2 + r \, N/2r$$

The Generalized Cost function for MFT-2 network is given by:

$$C = (8.2^{n+1} - N)$$

Where N = Number of Input and $n = \log_2 N$
The cost function of various network are as follows:

INDRA......................C= $4N(\log_2 N + 1)$

ASEN-2 ……......…. C= $3N (\log_2 N - 1)$

FT Network………..C= $(9.75 * 2^{n+1} - 54)$

The cost complexity of each network is analyzed below. The cost function is described as $\log_{10}C$ to convert them to smaller values. The cost of MFT-2 network is $(8.2^{n+1} - N)$. Fig. 11 shows the variation of $\log_{10}$(cost function) to $\log_2 N$. The graph clearly shows the MFT-2 network is more cost effective than other fault tolerant multistage

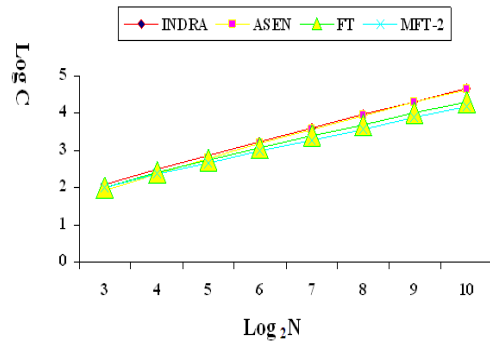interconnection network that's the advantage building large size networks.



**Fig.11: Comparison of $\log_{10}C$ to $\log_2 N$**

## 9    Conclusion

The MFT-2 network is a new cost effective fault tolerant irregular network with simple distributed control scheme. MFT-2 has alternate routing scheme in all stage except the last stage and has significant improvement in reliability and performance. Multiple path of variable length is available between same source and destination. Average path length is reduced. The MFT-2 network is more cost effective as compare to other fault tolerant networks.

### *References*

[1] Tse-yun Feng "A Survey of Interconnection Networks" IEEE Trans. On Computers Dec 1981 Vol 44, No.1 123-129.

[2] Daniel M. Dias and J. Robert Jump, "Analysis and Simulation of Buffered Delta Networks," *IEEE Trans. On Computers,* Vol.30, No. 4, pp. 273-282, April 1981.

[3] Chaun-LIN WU, "On a Class of Multistage Interconnection Networks", IEEE Trans. on Computers May 1980.

[4] Duncan H Lawrie, "Access and Alignment of Data in an Array Processor," IEEE Trans. on Computers,

Vol. 24, No. 12, pp. 1145-1155, December 1980.

[5]  Guofeng Hou, Tao Li, and Yulu Yang, "Research on Multistage Interconnection Network" *Proceedings of 2000 National Annual Conference on Computer Architecture by China Computer Federation,* Harbin, China, pp. 123-128, August 2000 ( in Chinese).

[6]  *Jehad Al-Sadi, Ahmad M. Awwad***,** "A New Fault-Tolerant Routing Algorithm for EOC Interconnection Network," WSEAS Transactions On Computers Issue 7, Volume 5, July 2006. pp. 1474-1480.

[7]  *D.C. Vasiliadis, G.E. Rizos, S.V. Margariti, L.E. Tsiantis,* **"**The Influence of Network Size on the Performance of Multistage Interconnection Networks," WSEAS Transactions on Communications, Issue 4, Volume 6, April 2007.pp 505 – 511.

[8]  *D.C. Vasiliadis, G.E. Rizos, S.V. Margariti, L.E. Tsiantis,* " Comparative Study of blocking mechanisms for Packet Switched Omega Networks, " Proceedings of the 6th WSEAS Int. Conf. on Electronics, Hardware, Wireless and Optical Communications, Corfu Island, Greece, February 16-17,2007,pp 18- 22.

[9]  Nitin, Vivek Kumar Sehgal, P.K. bansal, "On a MTTF Analysis of Fault Tolerant Hybrid MIN," WSEAS Transactions on Computer Research, Issue 2, Volume 2, February 2007.  pp. 130-138.

[10] P.K.Bansal, Kuldip Singh and R.C.Joshi, "Routing and path length algorithm for a cost effective four tree multistage interconnection network" Intl J. Elec 1992 vol 75 pp 107-112.

[11] P.K.Bansal, K.Singh, R.C.Joshi, G.P.Siroha, "Fault tolerant double tree network" Proc of Intl Conference IEEE INFOCOM April 1991 pp 462-468.

[12] K.Hwang, F.A.Brigges, "Computer Architecture and Parallel Processing" Tata McGraw Hill, Computer Organization and Architecture, 1984.

[13] R.Mittal, D.Cherion, P.J.Mohan, "Routing and Performance of the Double tree network" IEEE Proc. Computer Digit Tech., vol 142 no 2 March 1995.

[14] A.C.Aljundi, M.T.Kechadi, I.D.Scherson, "A study of an Evaluation methodology for unbuffered multistage interconnection networks" Proc IPDPS 2003.

[15] J.Sengupta, P.K. Bansal, "Performance of Regular and Irregular Dynamic MIN," proc of IEEE TENCON 1999.

[16] G. Adams, D.Agrawal and N.Siegel, "A Survey and comparison of fault tolerant multistage interconnection networks," IEEE computer pp 14-27 June 1987.

[17] Sandeep Sharma, P.K.Bansal, "A new fault tolerant multistage interconnection network" Proceeding of IEEE TENCON 2002.

[18] J.T.Blake and K.S.Trivedi, "Multistage interconnection Reliability" IEEE Transaction vol 38 no 11 Nov 1989.

[19] Feng T.Y,"A survey of MIN" IEEE Comp vol 14 no.12 Dec 1981 pp 743-758.

[20] Shen J.P, "Fault tolerance analysis of several interconnection networks" Proceedings of International Conference on Parallel Processing August 1982 pp 102-112.

[21] P.K.Bansal,Kuldip Singh, R.C.Joshi,G.P.Saroha, "Fault Tolerant Double Tree Multistage Interconnection Network", IEEE Transactions on Parallel and Distributed Systems, VOL. 24, NO. 13, JAN 1991.