

Article

TFBN: A Cost Effective High Performance Hierarchical Interconnection Network

M. M. Hafizur Rahman ^{1,*} , Mohammed Al-Naeem ¹ , Mohammed N. M. Ali ² and Abu Sufian ³ ¹ Department of Computer Networks & Communications, CCSIT, King Faisal University, Al Ahsa 31982, Saudi Arabia; naeem@kfu.edu.sa² KICT, International Islamic University, Malaysia (IIUM) Jalan Gombak, Kuala Lumpur 53100, Malaysia; moh.ali.exe@gmail.com³ Department of Computer Science, University of Gour Banga, Malda 732103, India; sufian@ugb.ac.in

* Correspondence: mhr Rahman@kfu.edu.sa

Received: 17 October 2020; Accepted: 18 November 2020; Published: 20 November 2020



Abstract: In order to fulfill the increasing demand for computation power to process a boundless data concurrently within a very short time or real-time in many areas such as IoT, AI, machine learning, smart grid, and big data analytics, we need exa-scale or zetta-scale computation in the near future. Thus, to have this level of computation, we need a massively parallel computer (MPC) system that shall consist of millions of nodes; and, for the interconnection of these massive numbers of nodes, conventional topologies are infeasible. Thus, a hierarchical interconnection network (HIN) is a rational way to connect huge nodes. Through this article, we are proposing a new HIN, which is a tori-connected flattened butterfly network (TFBN) for the next generation MPC system. Numerous basic modules are hierarchically interconnected as a toroidal connection, whereby the basic modules are flattened butterfly networks. We have studied the network architecture, static network performance, and static cost-effectiveness of the proposed TFBN in detail; and compared static network and cost-effectiveness performance of the TFBN to those of TTN, torus, TESH, and mesh networks. It is depicted that TFBN possesses low diameter and average distance, high arc connectivity, and temperate bisection width. It also has better cost-effectiveness and cost-performance trade-off factor compared to those of TTN, torus, TESH, and mesh networks. The only shortcoming is that the complexity of wiring of the TFBN is higher than that of those networks; this is because the basic module necessitates some extra short length link to form the flattened butterfly network. Therefore, TFBN is a high performance and cost-effective HIN, and it will be a good option for the next generation MPC system.

Keywords: massively parallel computer system; hierarchical interconnection network; TFBN; cost performance trade-off factor; static network performance; static cost effective factor; time cost effective factor; smart grid; energy storage system

1. Introduction

Computing power becomes a commodity like an internet in our modern daily life. Thus, like internet speed, the computation power is increasing, and like the demand for internet speed, the demand for computation power is also increasing. Researchers from all over the world are trying in various ways to mitigate this ever-increasing demand for computation power. These include multi-core and many-core, General-Purpose Graphics Processing Unit, grid computing, cloud computing, and so forth. However, these computing techniques are not sufficient to solve some modern life issues, viz. Internet of Thing (IoT), Artificial Intelligence (AI), machine learning, big data analytics [1]. In addition, these techniques are not enough to solve some computation-intensive problems and grand challenge

problems. For example, recently the impact of a planetary collision is revealed by the simulation, and this simulation is carried out by a COSmology MACHine (COSMA) supercomputer to have the simulation result in a reasonable time [2]. Definitely, a rational researcher could not expect a simulation will take a couple of months or even years. A reasonably short period of time can only be attained if the computer can operate so fast. For example, NVIDIA and the University of Florida plan to construct a new AI supercomputer that will deliver 700 petaflops by 2021 [3]. Thus, it is a clear indication that we need the next generation supercomputer called a massively parallel computer (MPC) system will compute in exaflops or even more like the zetta-flops level of performance [1,4–7].

Even for our life from disaster to modern society, we need massively parallel computation everywhere. For disaster prevention and mitigation, we need extensive computation. For healthcare, drug design, and personalized medicine, we need extensive computation too. For example, in order to stop the spreading of community transmission disease, COVID-19, we need to track the affected person and the spreading nature. To do this vital task in a timely manner, the enormous computing power of the MPC system along with artificial neural network (ANN), AI is being utilized [3,8]. In order to develop a safe and effective vaccine against COVID-19 or any other virus, the genome sequence of that virus needs to be analyzed. It is noteworthy that the analysis of the COVID-19 virus is quite challenging because of its genetic mutation in a new environment and geographical location. Indeed, there is also an MPC system that can play an inevitable role in the genome sequence analysis of COVID-19 [1,9].

In the power system, the IoT enabled smart grid is becoming more popular [10]. In smart grid topologies, the real-time monitoring system along with a high standard load forecasting model are very important [11]. For such models, massive computation power is necessary. Due to the advancement of renewable energy, the energy storage system (ESS) is also marking its importance with a great value. To exploit the benefits of ESS to the fullest extent, it is necessary to allow ESS to participate in the energy market where load forecasting plays a great deal of importance [12]. For such load forecasting, different single or hybrid predictive models have already been proposed, which requires high computation power to respond quickly. The MPC system can indeed play a significant role in such scenarios related to power system and energy.

Information and Communication Technology (ICT), both computation and communication, become the part and parcel of our daily life everywhere. We need computation ranging from mobile computing, cloud computing, to a massively parallel computer system and communication at any time and anywhere; it will especially be dominant in the post COVID-19 world. For example, the efficient use of irrigation water for agriculture in the desert by using sensors and IoT will be challenging and an MPC system can ease the total agriculture process in the desert. In addition, food security is quite important and challenging especially in a natural disaster and pandemic situation. It is worth noting that it is important for food importing countries like Saudi Arabia and other Gulf countries. The central management of food management ranging from crop production and reservation to distribution and flood protection will be convenient with the use of the MPC system. Gene analysis and modification can easily be done with the help of an MPC system for breeding new kinds of crops like rice, or wheat, or corn which can easily be grown in the desert with less water. The modeling and simulation of many complex and dynamic problems, generally which will be very expensive and sometimes impractical or even impossible to demonstrate physically, can easily be carried with the help of an MPC system.

The fastest supercomputer right now in the world is Fugaku yielded 513 petaflops performance, whereby 158,976 nodes are interconnected by 6D torus using Tofu D interconnect [1,13] followed by Summit and Sierra interconnected 4356 and 4320 nodes, respectively by fat-tree topology [14]. After that, Sunway Taihulight and Tianhe 2A interconnected 40,960 and 16,000 nodes, respectively, by mesh topology [15]. Here, it is to be noted that each node consists of many-core, i.e., each supercomputer or an MPC system uses millions of cores. For instance, the Fugaku uses 7,299,072 cores to construct this MPC system [13]. Thus, it is depicted that the contemporary supercomputer or an MPC system interconnects thousands of nodes using conventional topologies like

higher dimensional mesh or torus or fat-tree topology. A higher-dimensional regular torus network is a widely acceptable topology right now. The heat dissipation and thus the power consumption of a higher dimensional network is high compared to that of a lower-dimensional network. This phenomenon is revealed in the Fugaku supercomputer. Here, it requires 28.34 Megawatt power, which is 2.8 times more than its counter rival Summit [13].

As mentioned earlier, a low-dimensional network is better [16–18]. However, to interconnect millions of nodes, a conventional low-dimensional network is not suitable [19–23]. Therefore, low dimensional hierarchical interconnection networks (HIN) topology are indispensable to interconnect more than 1 million nodes for the future generation MPC system [24–26]. Considering many design attributes of an MPC system such as better inter-node communication performance (especially in terms of low latency and high throughput), low power consumption, and quick reconfiguration and high fault tolerance. The HIN combines the good attributes of many conventional networks. Many HINs have been proposed already viz., completely connected based HIN, hypercube based HIN, tree-based HIN, k -ary n -cube-based HIN, etc. However, none of these HINs are practically implemented even though they have some merits and demerits. Therefore, further investigation and exploration of HIN are needed.

In this research, we have proposed a time-cost effective and high performance and novel hierarchical interconnection network for the next generation MPC system. The proposed HIN consists of numerous basic modules, and these are connected among themselves in a hierarchical manner to form the consequent higher-level networks. The basic modules are a 2D-flattened butterfly network, whereas the higher-level network is a 2D-torus network. The higher-level networks are constructed by considering immediate lower-level networks as a network module and these modules are interconnected as a 2D-torus connection. The proposed HIN is the combination of flattened butterfly network and regular 2D-torus network, and thus it is called the Tori-connected Flattened Butterfly Network (TFBN). We have considered the 2D-torus network for higher-level networks because of its regular connection pattern, and a wrap-around link between end-to-end nodes provides an alternative path to reduce the congestion and contention and increase the throughput. Flattened butterfly networks used as a basic module also reduces the congestion and increases the throughput within the basic module because its many short length links increase the connectivity within the basic module and yield good performance [27,28]. The novelty and superiority of the proposed TFBN are assessed by static analysis of static network performance and static time-cost effective analysis.

The preliminary version of a study on TFBN is presented in a conference—some static network performance parameters of TFBN for the Level-2 only consisting of only 256 nodes [29]. The main objectives of this study are the versatile study of the static performance of networks by evaluating these static network performance parameters for the higher-level network, static cost-effectiveness analysis, and static cost-performance trade-off analysis. The practical acceptance of any interconnection network topology by the industry community depends on the detailed evaluation and analysis of static network performance, cost-effectiveness analysis, dynamic communication performance, fault tolerance performance, on-chip, and off-chip power consumption, computational intensive problem mapping in the network topology, and so forth. In this first step of detail study, we have considered static analysis of network performance, cost-effectiveness as well as cost-performance trade-off.

After the Introduction, the remainder of this article is organized as follows: Sections 2 and 3 discuss the network architecture of the proposed TFBN and the routing of the message within it, respectively. The detailed study of the static network performance of a TFBN is discussed in detail for both the lower level as well as a higher level in Section 4. The cost-effectiveness as well as the cost performance trade-off factor of this TFBN is analyzed in Section 5. Finally, the concluding remarks and the future works of this study are presented in Section 7.

of the upper-level networks of a TFBN demonstrated in Figure 2. $4 \times 2^q = 2^{(q+2)}$ free wires and its associated port in each BM used for higher-level connections. For vertical connection, we used $2(2^q)$ links and for the horizontal connection, we used $2(2^q)$, where inter-level connectivity represented by q , $q \in 0, 1, \dots, m$, when $q = 0$ represents the minimal inter-level connectivity, and, when $q = m$, it will be the maximum inter-level connectivity. In Figure 2, there is a (4×4) BM with $2^{(2+2)} = 16$ free ports, by considering $m = 2$. By considering $q = 0$, we will have four free ports for every upper-level connection divided into two for horizontal and two for vertical interconnections. Tied incoming and outgoing links will form a bidirectional link to connect two adjacent BMs of higher-levels of TFBN. Vertical in and vertical out wires tied to create vertical links and horizontal in and horizontal out wires connected to create horizontal links.

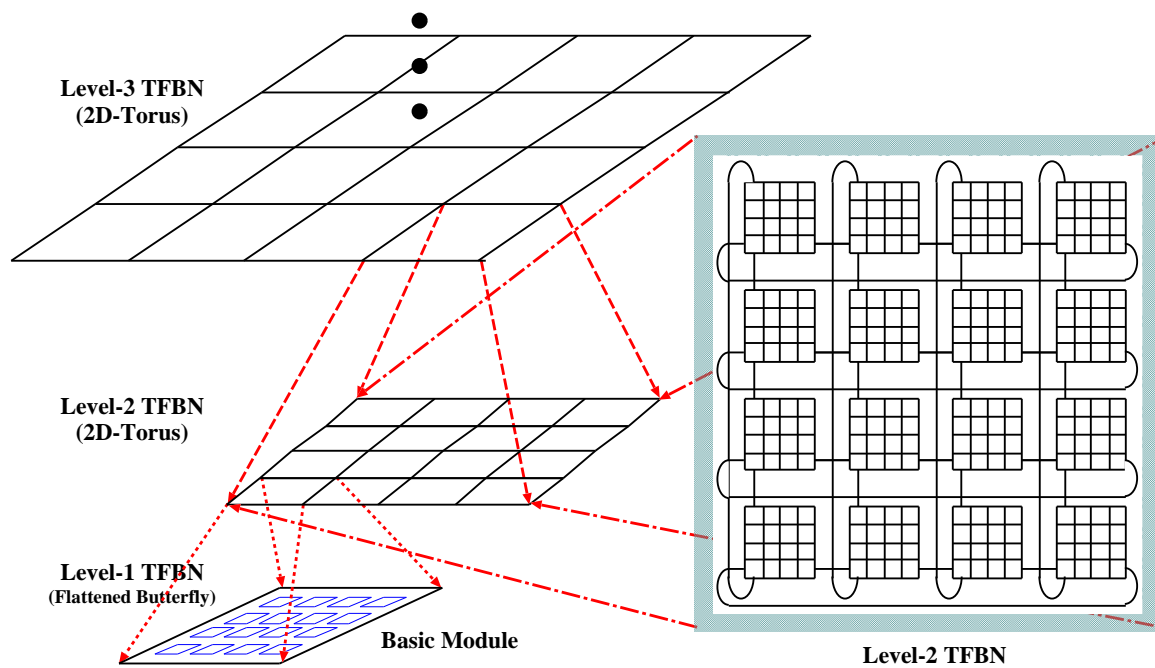


Figure 2. 4×4 higher level network of a TFBN.

TFBN represented by the value of m to decide the size of a BM, the level of hierarchy (L), and the inter-level connectivity (q) as TFBN(m, L, q). To guarantee better granularity, here we have considered $m = 2$. L_{max} is the maximum hierarchical level can be created from a $(2^m \times 2^m)$ BM and it is $L_{max} = 2^{(m-q)} + 1$. As an example, with $q = 0$ and $m = 2$, the highest level of hierarchy is $L_{max} = 2^{(2-0)} + 1 = 5$. This implies that we obtained Level-5 as a maximum possible level of TFBN by connecting (4×4) BM hierarchically. The number of nodes in each level (L) of a TFBN is $N = 2^{2mL}$. Therefore, the number of nodes can be obtained of connecting TFBN(m, L, q), where $L = 5$, $m = 2$, $q = 0$ is $N = 2^{2m(2^{(m-q)}+1)}$. Thus, 1,048,576 nodes could be interconnected together to create a massively parallel computer system based on using TFBN network topology.

3. Routing Algorithm for TFBN

The routing algorithm is vital to be defined in supercomputer systems to clarify the mechanism of sending and receiving messages between any two systems. Thus, it simplifies the path selection by a message inside the network to move from a source node to a destination node. The performance of the interconnection network, which is a vital mainstay in creating the massively parallel computer system that is affected widely by the proficiency of message routing and by the router design.

The routing procedure is divided into three phases to forward a message from its source to its destination in TFBN [30]. TFBN is a hierarchical interconnection network; therefore, forming any higher-level network will apply the same coordinate locations of the immediate lower level. In phase-1,

the message will be routed in the source BM, and it will be directed to the proper gate node to be sent to the higher-level network. In the second stage which is called Phase-2, the message will be routed through the higher-level network. Once it reaches the highest-level which leads to the destination BM, the routing operation will be reversed from the higher to the lower levels until the packet stabilizes in the final destination BM. In Phase-3, the message will be routed from the gate node in the destination BM to reach its final station in the destination node. The routing mechanism of Phase-1 and Phase-3 is almost the same. One is in the source BM and another one is in the destination BM.

A usual deterministic dimension-order routing protocol will be applied to ease the routing process in this network. In this protocol, the message will be routed first vertically and then it will be routed horizontally. In the source node, when the packet is generated, the destination will be checked. If the source, as well as the destination nodes, are in the same BM, then the procedure of routing will perform within the BM only. In this case, only Phase-1 will be applied. Otherwise, if the source and the destination nodes are in different BMs, the message shall be sent to a proper gate node that connects the BM to the level where the routing is carried out. In this situation, the three routing phases will be applied. The routing of messages by these three phases on TFBN is illustrated in the following Figure 3.

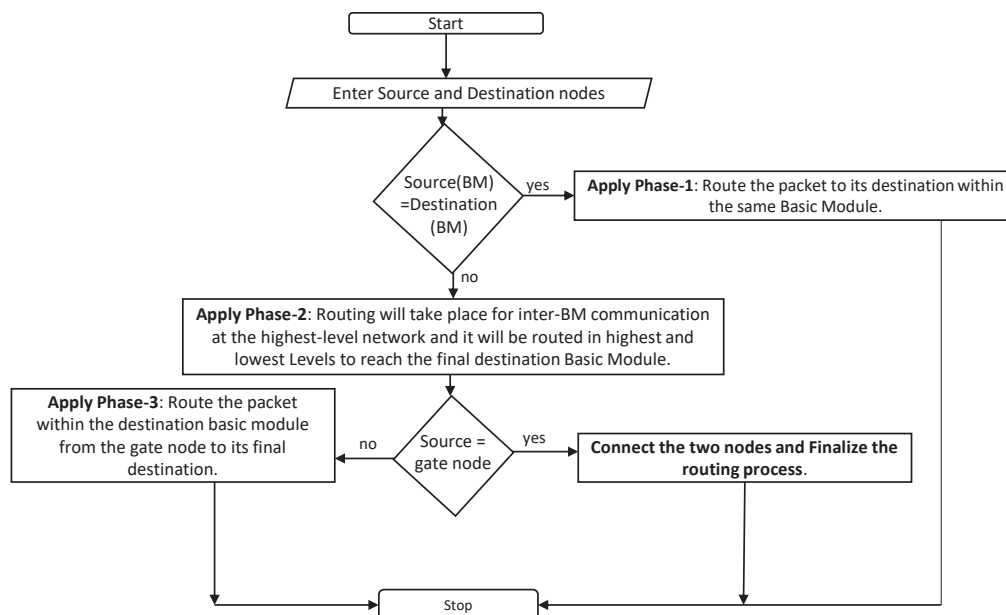


Figure 3. Illustration of message routing in a TFBN.

4. Static Network Performance Evaluation of a TFBN

In order to show the superiority of the proposed TFBN, we have compared its performance [31] with other interconnection networks. We have considered tori-connected torus network (TTN), tori-connected mesh network (TESH), torus, and mesh networks. Each and every network have their own structure and architecture. This is why it is quite difficult or even impossible to consider all the networks for performance comparison. This is why we have considered two similar hierarchical interconnection networks such as TTN [32,33] and TESH network [34,35], and two regular conventional networks for the performance comparison. Conventional mesh and torus networks are widely used for the practical implementation of the MPC system.

4.1. Node Degree

A node of an MPC system is consisting of a processing component, a router, and a memory. Many nodes are interconnected together using an interconnection network topology. The router of a node is used to connect to its neighboring node. The node degree or degree of a node is described as the number of links instigating from a node to connect all of its neighbor nodes. Different nodes of an interconnection network topology will have different degrees of the node. The maximum value of the degree of nodes of an interconnection network is called the node degree. Node degree is very important for the router of a node. A constant node degree is desirable for the good scalability of the massively parallel computer system. If the node degree is changing with the size of the network, then the scalability is difficult and expensive. Vice versa, it can be said that the scalability is not possible for the MPC system.

On the other hand, the cost of a router is directly proportional to the degree of a node. The lower node degree lowers the router cost. Usually, the number of the core of the processing element and the size of the memory of a node is constant in the chip-level design. Therefore, the degree of a node is directly proportional to the cost of that node. The lower the node, the lower the cost of a node. The degree of a node in different networks is tabulated in Table 1. It is illustrated that the node degree of proposed TFBN is constant and it is 8. The degree of a node the proposed TFBN is higher than that of mesh, torus, TESH, and TTN networks. However, this high cost of an individual node will result in a high performance of the proposed TFBN. It is shown in the next subsequent sub-section about the performance improvement and in the next section about the cost-performance effectiveness of the proposed hierarchical interconnection network TFBN.

Table 1. Comparison of static network performance of various networks.

	Node Degree	Diameter	Average Distance	Arc Connectivity	Bisection Width
256 Node					
2D-Mesh	4	30	10.67	2	16
2D-Torus	4	16	8	4	32
TESH	4	21	10.47	2	8
TTN	6	15	7.44	4	8
TFBN	8	10	5.75	4	8
4096 Node					
2D-Mesh	4	126	42.67	2	64
2D-Torus	4	64	32.00	4	128
TESH	4	32	17.80	2	8
TTN	6	24	12.60	4	8
TFBN	8	19	10.61	4	8

4.2. Diameter and Average Distance

Even though hop distance is not the actual distance between a node to node, it is a crucial parameter to evaluate the performance of any network in its graph-theoretic model. The maximum hop distance among all pairs of nodes in a network using the shortest-path algorithm is called the diameter of that network. The average hop distance by considering all the distinct pairs of nodes in a network using the shortest-path algorithm is known as average distance. The diameter signifies the upper bound of the latency or indicates the latency during the saturation throughput. The average distance signifies the latency at no load and the network throughput. The lower the diameter and average distance, the better the network in terms of dynamic performance communication. Therefore, the nature of the dynamic performance communication (low latency and high throughput) can be estimated well prior to the actual evaluation of it. If it is not acceptably good, then we can change some connection and evaluate these two parameters again.

Using simple dimension order routing, we have estimated the diameter and average distance of TFBN, TTN, and TESH network by simulation studies for both the Level-2 and Level-3 networks. On the other hand, we have calculated those parameters by their respective formula for the mesh and torus networks. The results are tabulated in Table 1. It is depicted that both the diameter and average distance of the TFBN are substantially low as compared to those of TTN, TESH, torus, and mesh networks. Particularly, it is revealed that the average distance of the TFBN is getting lower as the hierarchical level of the TFBN is increasing from Level-2 to Level-3.

4.3. Arc Connectivity and Bisection Width

The MPC system consisting of many nodes and they are connected by many links. It is quite normal that the node or the links may have a fault [36], or it may fail to operate properly. This fault or failure of links may cause the disjoint of the underlying interconnection network deployed to construct the MPC system. Two disjoint parts of the network might be just two parts or it might be two equal halves of the network. The minimal number of links that need to be cut or to be faulty to disconnect the network into two separate parts is known as arc connectivity. In addition, the minimum number of links that need to be cut or to be faulty to disconnect the network into two equal halves is called bisection width.

Arc connectivity is statically used to indicate the fault-tolerant capability of any interconnection network. The actual fault tolerance depends on the reconfiguration of the faulty nodes or links by redundant wires or links or by rerouting the packet by any alternative path. However, the ratio between arc connectivity and degree of node statically indicates the fault tolerance capability of any interconnection network. We have evaluated the arc connectivity of the proposed TFBN along with the TTN, TESH, torus, and mesh networks, and the evaluated result is charted in Table 1. It is portrayed that the arc connectivity of TFBN is equal to that of toroidal networks such as conventional torus and hierarchical TTN. In addition, it is double that of conventional mesh and hierarchical TESH network. It is noted from Table 1 that the high node degree of the TFBN results in a low ratio of arc connectivity and node degree. Therefore, static fault tolerance performance of the proposed TFBN is lower than conventional torus and TTN, and it is more than conventional mesh and hierarchical TESH networks.

The moderate value of bisection width is preferable. If it is low, then the merging of the yielded result of the computation outcome by the individual node will be slowed down because the bisection links will be congested due to heavy traffic flow between two halves of the interconnection network. This congestion will be severe for a bit complement traffic pattern where all the packets cross the bisection of the network. Thus, the bisection links will be heavily congested. Therefore, to reduce this congestion and contention, high bisection width is better. On the other hand, high bisection width will create the complexity in Very large-scale integration (VLSI) or wafer-stack implementation or the Network-on-chip (NoC) design of an interconnection network. Therefore, a kind of optimum and moderate value is expected for any interconnection network to be considered for the next generation MPC system that is neither too high nor too low. We have evaluated the bisection width of TFBN along with other networks considered in this paper. It is depicted in Table 1 that the TFBN has moderate bisection width like hierarchical TTN and TESH networks, and it is quite a bit lower than mesh and torus networks.

The question may arise as to why we have considered Level-2 (256 nodes) and Level-3 (4096 nodes) for the evaluation and comparison of performance. The proposed TFBN is the updated work of TTN, whereby the basic module of TTN is replaced by a flattened butterfly network. It is shown in our another study that TTN results in good performance for more than 1 million nodes compared to other networks [32,33]. For both the Level-2 and Level-3 networks, the proposed TFBN yielded superior performance to that of TTN. Therefore, it is believed that the TFBN will result in superior performance for Level 4 and 5 networks compared to other networks.

5. Static Cost Effectiveness Analysis of a TFBN

In this article, we have analyzed the proposed TFBN statically to show its superiority over that of other networks. In this section, we have described the studied the cost-effectiveness analysis [37] of the proposed TFBN to conduct further experiments on dynamic communication performance and the prototype hardware implementation. The parameters included for static cost-effectiveness analysis are static cost, wiring complexity, time-effective factor, cost-effective factor, and cost trade-off factors; these parameters are discussed below.

5.1. Cost

Even though the actual cost of an MPC system depends on its capital expenditure of nodes and their associated links to connect them, the product of the degree of node and network diameter is an acceptable static criterion to assess the cost of an MPC system. The degree of a node is directly proportionate to the cost of the router of a node, and the cost of a node is directly proportional to the cost of an MPC system. Connecting all the nodes by communication links forms the MPC system. The cost of an MPC system inculcating by node degree is to connect the communication links between the nodes to form the MPC system. The diameter of a network is the highest hop distance from the source to the destination using the shortest path algorithm among all distinct source–destination pairs. Therefore, the product of these two parameters named degree and diameter is a wise criterion to compare different interconnection network topologies.

The cost of different networks with two different sizes are estimated and charted in Table 2. It is presented that the cost of TFBN is lower than all Level-2 hierarchical networks and conventional mesh and torus networks with 256 nodes. However, the TFBN is a little bit costlier than Level-3 hierarchical networks, and much less than conventional mesh and torus networks with 4096 nodes. This slightly high cost incurred due to the high degree of node even though the diameter is quite lower than those hierarchical networks considered in this article.

Table 2. Comparison of static cost, wiring complexity, and cost performance trade-off factor for various interconnection networks.

	Cost	Wiring Complexity	Cost Performance Trade-off Factor
256 Node			
2D-Mesh	120	480	0.25
2D-Torus	64	512	0.50
TESH	84	416	0.31
TTN	90	544	0.85
TFBN	80	800	2.50
4096 Node			
2D-Mesh	504	8064	0.0625
2D-Torus	256	8192	0.1250
TESH	128	6680	0.2039
TTN	144	8736	0.5332
TFBN	152	12832	1.3191

5.2. Wiring Complexity

The wiring complexity is represented as the number of communication links required for interconnecting entire nodes of an MPC system. It is a static cost estimation for the communication links cost. The communication link cost depends on the length of the links and their underlying interconnection ranging from the VLSI chip level to the board level, cabinet-level, and system level. The chip level designed is ignored here because we have considered our design using node. Therefore, our consideration is from the board level to the system level. In static wiring complexity evaluation,

we have considered the number of communication links and not the length of the links. Thus, both the long length link and short length link are considered as an equal cost.

The proposed TFBN consists of numerous basic or primary modules whereby this primary module is a 2D-torus network. The wraps links in these basic modules result in a bit high number of communication links, which is, as mentioned before, the wiring complexity. The wiring complexity of a TFBN is $\left[\# \text{ of links in a BM} \times k^{2(L-1)} + \sum_{x=2}^L 2(2^x) \times k^{2(L-1)} \right]$, where L is the level number. As tabulated in Table 2, the wiring complexity of the proposed TFBN is more than that of mesh, torus, TESH, and TTN networks. These extra-short length wrap-around links in the basic module result in short diameter and low average distance.

5.3. Static Cost Performance Trade-off Analysis

Before going to prototype implementation or practical implementation of a huge expensive MPC system, many early stages of evaluation are crucially needed to investigate the suitability of the proposed interconnection network. In this study, we have studied the very first stage of evaluation and statically evaluated the suitability of the proposed TFBN using many static parameters. We have considered many static parameters and the last one is the static cost-performance trade-off factor (CPTF). It is mentioned earlier that the cost incurring parameters of an MPC system is a router in the node (i.e., node degree) and a number of communications links and their total length to layout the system design (i.e., wiring complexity). The most decisive parameter for the static network performance is the diameter; this indicates the upper bound of the communication latency. The diameter again depends on the network size, i.e., the total number of nodes. Considering all these parameters, the CPTF of an interconnection network topology is expressed by the following Equation (1). The higher the outcome of Equation (1), the more suitable the interconnection network:

$$CPTF = \frac{\text{Node Degree} \times \text{Wiring Complexity}}{\text{Diameter} \times \text{Total \# of Nodes}} \quad (1)$$

We have assessed the factor of static cost-performance trade-off of TFBN along with TTN, TESH network, 2D-torus, and 2D-mesh networks, and the results are tabulated in Table 2 for both 256 nodes and 4096 nodes network. It is stated that the CPTF of a TFBN network is far higher than that of Torus, Mesh, TESH, and TTN networks for both sized network (256 nodes and 4096 nodes).

5.4. Cost Effective Factor

A massively parallel computer system is constructed by connecting numerous nodes interconnected using communication links under an interconnection network topology. The major cost of a node is incurred because of its processing elements, router, and memory. Considering the price of communication links or cables in creating an MPC system is also one of the important factors to calculate the actual cost of an MPC system. The price of the wires has a crucial influence on the cost of building MPC systems; therefore, finding a new parameter using the cost of these wires to evaluate these systems is important.

Usually, the parallel computer system is used to reduce the processing time and speed up the operation. The evaluation of the performance of MPC systems takes into account speedup and efficiency to assess these systems. The more nodes that are interconnected together, the more speed that will be attained. In addition, surely the ratio is not linear. The speedup and thus the efficiency of an MPC system depend on the proficient use of the communication links used to interconnect the MPC system.

The cost-effective factor (CEF) is considered as a crucial static parameter to assess the feasibility of an MPC system. This parameter considers the cables count as a function of the processors' count. Therefore, it denotes $CEF(p, l)$, where p represents the number of processing elements or cores and l represents the number of links. The proportion of wires count to nodes count denoted as $G(p) = (\text{Total \# of links})/(\text{Total \# of nodes})$. The cost of processing elements is represented as C_p , and the cost of communication links is represented as C_l . The parameter ρ is outlined as the ration between C_l and C_p , and $0 \leq \rho \leq 1$. To evaluate the CEF parameter, we have considered the homogeneous parallel computer systems whereby the cost of processors and cables is maintained as a similar ration. Considering p , l , $G(p)$, C_p , and C_l , and ρ , the CEF of an interconnection network topology is expressed as the following Equation (2):

$$CEF = \frac{1}{1 + \rho \times G(p)} \quad (2)$$

We have assessed time cost-effective factors of TFBN along with the TTN and TESH network, 2D-torus, and 2D-mesh networks for different ρ ranging from 0.1 to 1.0 and plotted in Figure 4. We have considered two different size networks such as 256 nodes and 4096 nodes. Figure 4a portrayed the CEF of different networks for 256 nodes and Figure 4b portrayed the same for 4096 nodes. For all the values of ρ , ($0.1 \leq \rho \leq 1.0$), the CEF of TFBN is substantially less than that of TTN, TESH, torus, and mesh networks. From $\rho = 0.3$ and onward, the CEF of the TFBN is significantly lower than those networks for both 256 node and 4096 node networks. This proves that TFBN is a profitable network compared to these networks. This qualifies TFBN to be a good choice compared to many conventional and hierarchical interconnection networks to frame future generation supercomputers or MPC systems.

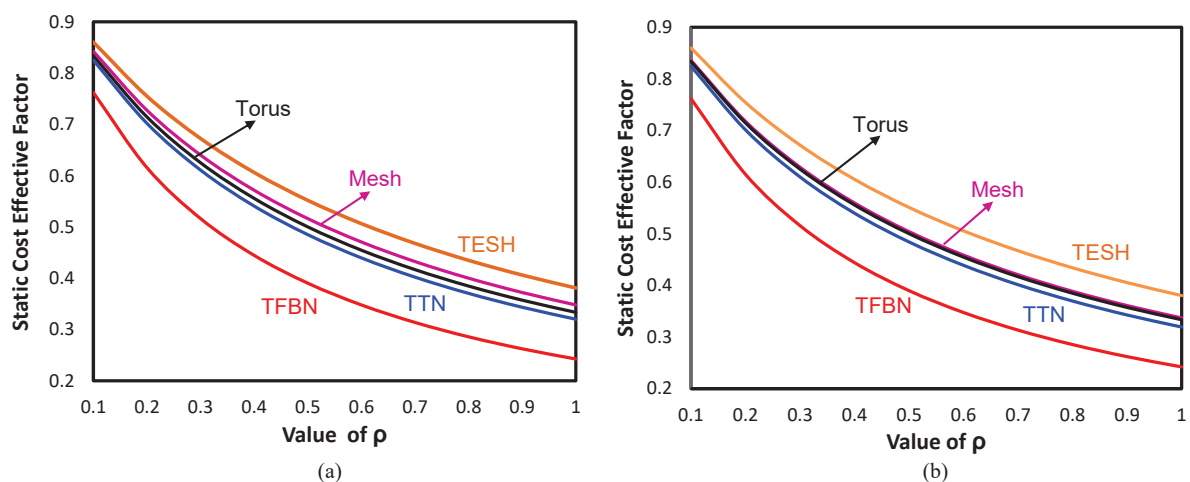


Figure 4. Comparison of static cost-effective factor of various interconnection networks considering different values of ρ (a) network size 256 nodes; (b) network size 4096 nodes.

5.5. Time Cost Effective Factor

The main purpose of the use of an MPC system is to solve the computational intensive problem or a grand challenge problem in a reasonably short period of time. This time depends on four factors, viz., the processing speed of an individual node, the dividing of the tasks among the nodes an MPC system and coordination of these tasks between nodes, the underlying routing algorithm to transfer the flits between source and destination nodes, and after the execution of each individual node task merging these individual execution outcomes to conclude with the final result. The fast time implies that the MPC system fast individual node and the underlying interconnection network has good performance in terms of less blocking, delay, and congestion. Thus, time is also an important factor in measuring the efficiency of an MPC system. Cost-effective factors considered along with time factors to determine

the time cost-effective factor (TCEF) of an interconnection network. Therefore, the TCEF is a very important and essential parameter to statically assess the feasibility of an MPC system, especially its underlying interconnection network topology. TCEF of an interconnection network is expressed by Equation (3):

$$TCEF(p, T_p) = \frac{1 + \sigma T_1^{\alpha-1}}{1 + \rho G(p) + \frac{\sigma}{p} T_p^{\alpha-1}} \quad (3)$$

TCEF of a network denotes as $TCEF(p, T_p)$ where p denotes the number of processor in an MPC system, and T_p denotes time taken to solve a problem using that MPC system. T_1 denotes the time taken to solve the same problem using a single processing element. Here, $G(p) = \frac{(\text{Total \# of links})}{(\text{Total \# of nodes})}$ and $\rho = \frac{C_l}{C_p}$, whereby $0 \leq \rho \leq 1$. The values of α and σ are constants and considered equal to 1 for considering linear time penalty in T_p . Including all of these values, Equation (3) is simplified as below:

$$TCEF(p, T_p) = \frac{2}{1 + \rho G(p) + \frac{1}{p}} \quad (4)$$

We have assessed the time cost-effective factor of TFBN along with TTN, TESH network, 2D-torus, and 2D-mesh networks for different ρ ranging from 0.1 to 1.0 and plotted in Figure 5. Like CEF, here, we have also considered two different size networks such as 256 nodes and 4096 nodes. Figure 5a depicted the TCEF of different networks for 256 nodes and Figure 5b depicted the same for 4096 nodes. For all the values of ρ , ($0.1 \leq \rho \leq 1.0$), the TCEF of TFBN is significantly lower than that of TTN, TESH, torus, and mesh networks. From $\rho = 0.3$) and onward, the TCEF of the TFBN is expressively lower than those networks for both 256 node and 4096 node networks. This proves that TFBN is a promising hierarchical interconnection network as compared to these networks. This qualifies TFBN to be a plausible alternative choice compared to many conventional and hierarchical interconnection networks to build next-generation supercomputers or MPC systems.

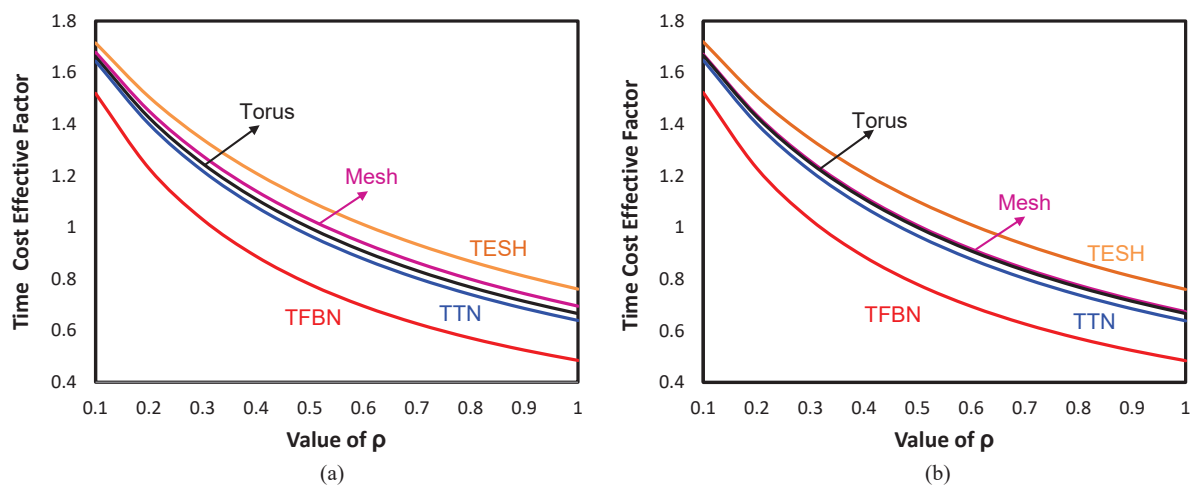


Figure 5. Comparison of static time cost-effective factors of various interconnection networks considering different values of ρ (a) network size 256 nodes; (b) network size 4096 nodes.

6. Some Generalization

Many interesting applications and grand challenging problems are discussed in the Introduction. To solve those problems, we need an MPC system. In addition, the success of an MPC system heavily depends on the reliable and suitable interconnection network. The reliability and suitability of an interconnection network topology are analyzed and assessed from different points of views for the consideration of its pragmatic realization. The very first step of the evaluation of an interconnection

network is the assessment of its graph-theoretic properties known as static network performance. These static network performance parameters revealed whether the proposed interconnection network will result in a good performance or not in other aspects.

For example, node degree indicates the cost of a router used in a node; and each node requires a router. Therefore, the cost of an MPC system depends on the node degree and wiring complexity. The low value of node degree and less wiring complexity clearly indicate the lower manufacturing cost of an MPC system. A low value in the distance parameter such as diameter and average distance indicates the possibility to have good dynamic communication performance. The lower hop distance parameter (low diameter and average distance) results in low latency and high throughput. The actual fault tolerance of an MPC system depends on the reconfiguration and routing with redundant resources. However, the ratio between the arc connectivity and node degree results in the static fault tolerance performance. The higher the ratio, the better the fault-tolerance of that network. The proposed TFBN results in significantly low diameter and average distance and reasonably better fault tolerance with the detriment of high node degree and wiring complexity as compared to other networks considered in this paper. However, the static cost (as portrayed in Table 2) of the TFBN is lower than that of other networks.

The feasibility and cost-effectiveness analysis are also imperative for the consideration of any new interconnection network topology. Static assessment and evaluation of cost-effectiveness along with the cost-performance trade-off factor without any capital expenditure are the good criteria to compare and contrast any new interconnection network with other contemporary networks [38]. The better performance in this static evaluation will instigate the next level of investigation for the suitability of the proposed network. As depicted in this paper, the cost-effective factor and time-cost effective factor of the proposed TFBN are lower than all of the networks, especially $\rho = 0.3$ and onward. In addition, the cost performance trade-off factor of TFBN is appreciably higher than all the interconnection networks.

The initial investigation instigates us for further exploration of the proposed TFBN; and the static hop distance parameters signify that the dynamic communication performance will also be better in terms of low message latency and high network throughput. However, it seems that TFBN will reveal good performance and will be a good choice for the next generation MPC system to overcome the challenging problem. The only shortcoming is the high wiring complexity. Many short length links in the basic module and their layout in the chip will incur a bit more power consumption. The chip-level layout of these links and the reduction of the capacitive effect between two links can reduce this extra power consumption.

The main concept and contribution in the paper is a new HIN called TFBN and its static network performance and cost-effective performance evaluation. With respect to static network performance and cost effective performance evaluation and analysis, the proposed TFB is superior to other conventional and hierarchical interconnection networks. The open issues and challenges have been mentioned earlier.

The success of an MPC system is how fast and efficiently it can execute and process a complex task, and the time required for this purpose depends on the execution time and communication and coordination time among the nodes. The execution time is constant, and it depends on the number of cores in the processor and its clock cycles. The execution time is fixed for a particular node of an MPC system. Communication time depends on the number of steps for communication and coordination. Therefore, the total required time is proportional to the number of steps required for communication and coordination.

The practical applicability of an interconnection network topology is usually justified by assessing the number of communication and coordination steps required for the benchmark computational intensive problem, viz., fast Fourier transform, solution of a partial differential equation, bitonic merge, finding the maximum, etc. It is believed that the proposed TFBN will result in less communication steps to solve these computational intensive problems because of the low diameter and average distance. A similar phenomenon is observed in our another study [39].

The implementation of a massively parallel computer system is the integration of many chips in the chip level integration to create a node. The interconnection of many nodes in the board level interconnection to make a board, interconnection of many boards in the cabinet level interconnection is to make a cabinet, and finally interconnection of many cabinets in the system level interconnection is to build the MPC system. In the TFBN, Level 1 (BM) is considered as the board, Level-2 is considered as the cabinet, and Levels 3, 4, and 5 are considered as system level interconnection. As mentioned before, Level-4 is interconnected using Level-3 as a sub-net module and similarly Level-5 is interconnected using Level-4 as a sub-net module. In this research, we have considered the static network performance and static cost effective analysis of the proposed TFBN. This is the very first stage of analysis, and, in the next step of this research, we will evaluate the dynamic communication performance by a flit-level simulator.

7. Conclusions

In this article, we have statically proposed and studied details of the TFBN. The architectural structure of the TFBN has been discussed in detail, and its superiority is depicted by evaluating many static parameters and comparing these parameters with Torus, Mesh, TESH, and TTN networks. The static parameters considered in this paper to show the preeminence of the TFBN are hop distance parameters such as diameter and average distance; connectivity parameters such as node degree, arc connectivity, and bisection width; cost parameter such as cost and wiring complexity; and cost-effectiveness parameters such as cost-effective factor (CEF), time cost-effective factor (TCEF), and cost-performance trade-off factor (CPTF).

Results evaluated of the above parameters revealed that TFBN possesses several attractive features. These include quite a low diameter and average distance, high static fault-tolerant and moderate bisection width, substantially low CEF and TCEF, and significantly high CPTF compared to those of TTN, TESH, torus, and mesh networks. These benefits are attained with the cost of a high node degree and wiring complexity. The flattened butterfly network as a basic module needs a few more short length links which result in the complexity of the wiring of the proposed TFBN being a bit high as compared to other networks considered in this paper. Using these extra links for the interconnection of the basic module in turns increases the node degree. It is shown that this high node degree increased the static cost of the proposed TFBN marginally higher than TTN and TESH networks. However, it is still far lower than the torus and mesh networks. From the cost-performance trade-off analysis, it is divulged that the TFBN is highly cost-effective compared to other networks. Therefore, TFBN will be a promising HIN topology to construct an MPC system that supports exa-scale or zetta-scale computation power.

The proposed TFBN was statically assessed and the performance was analyzed from various points of view. This is the first-stage research on TFBN. There is a long way to go for the consideration of TFBN by the industry community. Even though the diameter and average distance indicate the yielding of high throughput and low latency, evaluation of these two parameters (latency and throughput) has not yet been evaluated yet. Evaluation of these two parameters using a flit level simulator under deterministic dimension order routing and their improvement using adaptive routing algorithms are kept as the immediate next future step of evaluation [40]. Statically, the TFBN is cost-effective, and its cost performance trade-off is also quite good; however, prototype implementation by FPGA is necessary to assess the actual cost of the proposed TFBN [41]. The study can be implemented in various application domains such as the IoT enabled smart grid system and evolving ESS technologies.

Author Contributions: Conceptualization, M.M.H.R.; methodology, M.M.H.R. and M.A.-N.; software, M.M.H.R.; validation and formal analysis, M.M.H.R., M.A.-N., and A.S.; writing—original draft preparation, M.M.H.R. and M.A.-N.; writing—review and editing, M.M.H.R., M.N.M.A., and A.S.; visualization, M.A.-N. and A.S.; supervision, funding acquisition, project administration, M.M.H.R. and M.N.M.A. All authors have read and agreed to the published version of the manuscript.

Funding: The authors extend their appreciation to Deputyship for Research & Innovation, Ministry of Education in Saudi Arabia for funding this research work through the project number IFT20046.

Acknowledgments: Some static network performance analysis for the Level-2 network is presented in the 2nd international conference EICT, Khulna, Bangladesh, 2015. This paper focuses on the more detailed study of static network performance for the upper-level TFBN along with the cost-effectiveness analysis. The authors extend their appreciation to Deputyship for Research & Innovation, Ministry of Education in Saudi Arabia for funding this research work through the project number IFT20046.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Dongarra, J. *Reports on the Fujitsu Fugaku Fugaku System*; Tech Report No. ICL-UT-20-06; University of Tennessee, Oak Ridge National Laboratory: Knoxville, TN, USA, 2020.
2. Kegerreis, J.A.; Eke, V.R.; Massey, R.J.; Teodoro, L.F.A. Atmospheric Erosion by Giant Impacts onto Terrestrial Planets. *Astrophys. J. Am. Astron. Soc.* **2020**, *897*, 161. [CrossRef]
3. University of Florida, Nvidia Plan Fastest AI Supercomputer in Academia. Available online: <https://insidehpc.com/2020/07/> (accessed on 23 July 2020).
4. Wang, B. Ten Exascale Supercomputers by 2023. Available online: <https://www.nextbigfuture.com/2018/10/ten-exascale-supercomputers-by-2023.html> (accessed on 23 July 2020).
5. Dongarra, J.; Gottlieb, S.; Kramer, W.T. Race to Exascale. *Comput. Sci. Eng.* **2019**, *21*, 4–5. [CrossRef]
6. Beckman, P. Looking Toward Exascale Computing. In Proceedings of the Keynote Speech, International Conference on Parallel and Distributed Computing, Applications and Technologies (PDCAT'08), Dunedin, New Zealand, 1–4 December 2008.
7. Nagel, W.E. From TERA- to PETA- to EXA-Scale Computing: What does that mean for our Community? In Proceedings of the Keynote Speech in the 10th IASTED PDCN, Innsbruck, Austria, 15–17 February 2011.
8. Car, Z.; Segota, B.; Andelic, N.; Lorencin, I.; Mrzljak, V. Modeling the Spread of COVID-19 Infection Using a Multilayer Perceptron. *Hindawi Comput. Math. Methods Med.* **2020**, *2020*, 1–10. [CrossRef] [PubMed]
9. Shaheen, I.I. Open to Serve COVID-19 Research around the Kingdom of Saudi Arabia. Available online: <https://www.kaust.edu.sa/en/news/shaheen-ii-open-to-serve-covid-19-research-around-the-kingdom> (accessed on 23 July 2020).
10. Sakib, N.; Hossain, E.; Ahamed, S.I. A Qualitative Study on the United States Internet of Energy: A Step Towards Computational Sustainability. *IEEE Access* **2020**, *8*, 69003–69037. [CrossRef]
11. Al Mamun, A.; Sohel, M.; Mohammad, N.; Sunny, M.S.H.; Dipta, D.R.; Hossain, E. A Comprehensive Review of the Load Forecasting Techniques Using Single and Hybrid Predictive Models. *IEEE Access* **2020**, *8*, 134911–134939. [CrossRef]
12. Hossain, E.; Faruque, H.M.R.; Sunny, M.; Haque, S.; Mohammad, N.; Nawar, N. A Comprehensive Review on Energy Storage Systems: Types, Comparison, Current Scenario, Applications, Barriers, and Potential Solutions, Policies, and Future Prospects. *Energies* **2020**, *13*, 3651. [CrossRef]
13. Supercomputer Fugaku. Available online: <https://www.top500.org/lists/top500/2020/06/> (accessed on 23 July 2020).
14. Oak Ridge National Laboratory. Available online: <https://newatlas.com/oak-ridge-most-powerful-supercomputer-summit/54982/> (accessed on 23 July 2020).
15. Fu, H.; Liao, J.; Yang, J.; Wang, L.; Song, Z.; Huang, X.; Yang, C.; Xue, W.; Liu, F.; Qiao, F.; et al. The Sunway TaihuLight Supercomputer. System and Applications. *Sci. China Inf. Sci.* **2016**, *59*, 072001. [CrossRef]
16. Alam, M.; Varshney, A.K. A Comparative Study of Interconnection Network. *Int. J. Comput. Appl.* **2015**, *127*, 37–43. [CrossRef]

17. Moudi, M.; Othman, M. On the relation between network throughput and delay curves. *Automatika* **2020**, *61*, 415–424. [\[CrossRef\]](#)
18. Moudi, M.; Othman, M.; Lun, K.Y.; Rahiman, A.R.A. x-Folded TM: An efficient topology for interconnection networks. *J. Netw. Comput. Appl.* **2016**, *73*, 27–34. [\[CrossRef\]](#)
19. Prasad, N.; Mukherjee, P.; Chattopadhyay, S.; Chakrabarti, I. Design and evaluation of ZMesh topology for on-chip interconnection networks. *J. Parallel Distrib. Comput.* **2018**, *113*, 17–36. [\[CrossRef\]](#)
20. Camarero, C.; Martinez, C.; Beivide, R. L-Networks: A Topological Model for Regular 2D Interconnection Networks. *IEEE Trans. Comput.* **2013**, *67*, 1362–1375. [\[CrossRef\]](#)
21. Andujar, F.J.; Villar, J.A.; Sanchez, J.L.; Alfaro, F.J.; Duato, J. N-Dimensional Twin Torus Topology. *IEEE Trans. Comput.* **2015**, *64*, 2847–2861. [\[CrossRef\]](#)
22. Seo, J.H.; Sim, H.; Park, D.H.; Park, J.W.; Lee, Y.S. One-to-One Embedding between Honeycomb Mesh and Petersen-Torus Networks. *Sensors* **2011**, *11*, 1959–1971. [\[CrossRef\]](#) [\[PubMed\]](#)
23. Mnejja, S.; Aydi, Y.; Abid, M.; Monteleone, S.; Catania, V.; Palesi, M.; Patti, D. Delta Multi-Stage Interconnection Networks for Scalable Wireless On-Chip Communication. *Electronics* **2020**, *9*, 913. [\[CrossRef\]](#)
24. Faisal, F.A.; Rahman, M.M.H.; Inoguchi, Y. A new power efficient high performance interconnection network for many-core processors. *J. Parallel Distrib. Comput.* **2017**, *101*, 92–102. [\[CrossRef\]](#)
25. Ali, M.N.M.; Rahman, M.M.H.; Nor, R.M.; Behera, D.K.; Sembok, T.M.T.; Miura, Y.; Inoguchi, Y. SCCN: A Time Effective Hierarchical Interconnection Network for Network-on-Chip. *Mob.Netw. Appl.* **2019**, *24*, 1255–1264. [\[CrossRef\]](#)
26. Rahman, M.M.H.; Shah, A.; Fukushi, M.; Inoguchi, Y. HTM: A New Hierarchical Interconnection Network for Future Generation Parallel Computers. *IETE Tech. Rev.* **2016**, *33*, 93–104. [\[CrossRef\]](#)
27. Kim, J.; Dally, W.J.; Abts, D. Flattened butterfly: A cost-efficient topology for high-radix networks. *IEEE J. Solid-State Circuits* **2007**, *43*, 29–41.
28. Kim, J.; Dally, W.J.; Abts, D. Flattened Butterfly: A Cost-Efficient Topology for High-Radix Networks. In Proceedings of the 34th Annual International Symposium on Computer Architecture (ISCA), San Diego, CA, USA, 9–13 June 2007; pp. 126–137.
29. Sohaini, M.H.; Rahman, M.M.H.; Nor, R.M.; Sembok, T.M.T.; Akhand, M.A.H.; Inoguchi, Y. A Low Hop Distance Hierarchical Interconnection Network. In Proceedings of the 2nd International Conference on Electrical Information and Communication Technologies (EICT), Khulna, Bangladesh, 10–12 December 2015; pp. 43–47.
30. Holmark, R.; Kumar, S.; Palesi, M.; Mekia, A. HiRA: A Methodology for Deadlock Free Routing in Hierarchical Networks on Chip. In Proceedings of the 3rd ACM/IEEE NOCS, San Diego, CA, USA, 10–13 May 2009; pp. 2–11.
31. Abd-El-Barr, M.; Al-Somani, T.F. Topological Properties of Hierarchical Interconnection Networks: A Review and Comparison. *J. Electr. Comput. Eng.* **2011**, *2011*, 189434. [\[CrossRef\]](#)
32. Rahman, M.M.H.; Inoguchi, Y.; Sato, Y.; Horiguchi, S. TTN: A High Performance Hierarchical Interconnection Network for Massively Parallel Computers. *IEICE Trans. Inf. Syst.* **2009**, *E92D*, 1062–1078. [\[CrossRef\]](#)
33. Rahman, M.M.H.; Sato, Y.; Inoguchi, Y. High and stable performance under adverse traffic patterns of tori-connected torus network. *Comput. Electr. Eng.* **2013**, *39*, 973–983. [\[CrossRef\]](#)
34. Jain, V.K.; Ghirmai, T.; Horiguchi, S. TESH: A new hierarchical interconnection network for massively parallel computing. *IEICE Trans. Inf. Syst.* **1997**, *E80-D*, 837–846.
35. Miura, Y.; Kaneko, M.; Rahman, M.M.H.; Watanabe, S. Adaptive Routing Algorithms and Implementation for TESH Network. *Commun. Netw.* **2013**, *5*, 16. [\[CrossRef\]](#)
36. Bossard, A.; Kaneko, K. Cluster-Fault Tolerant Routing in a Torus. *Sensors* **2020**, *20*, 3286. [\[CrossRef\]](#)
37. Kim, J.; Dally, W.J.; Scott, S.; Abts, D. Cost-efficient dragonfly topology for large-scale systems. *IEEE Micro* **2009**, *29*, 33–40. [\[CrossRef\]](#)
38. Yunus, N.A.M.; Othman, M.; Hanapi, Z.M.; Lun, K.Y. Reliability Review of Interconnection Networks. *IETE Tech. Rev.* **2016**, *33*, 596–606. [\[CrossRef\]](#)
39. Rahman, M.M.H.; Horiguchi, S. A Deadlock-Free Routing Algorithm using Minimum Number of Virtual Channels and Application Mappings for Hierarchical Torus Network. *Int. J. High Perform. Comput. Netw.* **2006**, *4*, 174–187. [\[CrossRef\]](#)

40. Rahman, M.M.H.; Sato, Y.; Inoguchi, Y. High Performance Hierarchical Torus Network under Adverse Traffic Patterns. *J. Netw.* **2012**, *7*, 456–467. [[CrossRef](#)]
41. Fukase, N.; Miura, Y.; Watanabe, S.; Rahman, M.M.H. The Performance Evaluation of a 3D Torus Network using Partial Link-Sharing Method in NoC Router Buffer. *IEICE Trans. Int. Syst.* **2017**, *E100-D*, 2478–2498. [[CrossRef](#)]

Publisher’s Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



© 2020 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).