NOC ROUTER AND TOPOLOGIES

The chapter explains the 2D and 3D router functionality and its use in the configuration in NoC design. The chapter also explains the detailed description of Mesh, Torus and Ring topology.

3.1 Router Architecture (2D)

The router is the device used to forward the data packets in communication network. The data packets are transferred form one source router to another destination router in the internetwork. Router connected to two or many data lines inside the network to transfer the packets to destination. When the data packets are arrived on these lines the router read the address of then destination network to deliver the data packets. In the routers can transmit the packets from one network to another network in this way. The common used routers are small office routers and homes which are passing simple IP based routers between the internet and home computers. The example of a router can be DSL or owner's cable that provides the connectivity to the internet through internet service provider (ISP). Some examples of the more safe and sophisticated routers are enterprise routers which are helpful to connect large ISP networks and large business networks as power main routers that can forward the data at very high speed rate in the direction of the optical fiber lines along with internet backbone. The routers are designed based on dedicated routing and these are dedicated hardware device. There is the existence of software based routers also.

NoC provides the architecture and communication model in which multiple nodes can communicate and utilize the resources. Due to this reason, it is possible to design and implement the chip working under independent resource utilization in one block and multiple blocks as the working node or processing element of the communicating network. The flexibility and scalability in the chip design is cable to provide the feasibility to configure the NoC with different cluster size and workloads. The NoC architecture consists of the multiple resources which are designed and arranged in particular topology. The common and important topology used for the NoC is mesh topology. In mesh topology the local interconnections are existing between the switches and resources. The interconnections are not depending on the communicating network cluster size. The routing of the 2D network can be done in easiest way and gives the overall scalability, large bandwidth and short period clock cycle. The NoC have the switches are resources which are interconnected in such a manner that they can communicate directly among each other using packet data transfer technique. The resource is the computation and storage unit. The NoC switch is used to provide the path or routing to incoming traffic and buffered the same traffic between the other resources. There are many inputs and output channels or interconnections through which the switch can connect to other neighbour's switches. Each interconnection channel has two or more buses in one direction between two switches as point to point connections.

The incoming traffic is handled using wait and go technique for queue and controlled by switch. The size of mesh network depends on the clusters and their associated memory. The fabrication technology can change chip size. The size of chip depends on the resources utilization and number of hardware resources. The applications and system volume demand the larger bandwidth but it is the challenge to the design engineer to provide the chip design and related communication in the allocated or associated bandwidth.



Figure 3.2 Crossbar switch of 2D NoC router



All the processing elements and the routers in NoC architecture are connected with the help of several connecting wires. The processing elements are processing the data under flow control technique with fixed length. The 2D router is shown in [Figure 3.1]. It can accept the data for five distinguish ports as the part of core: east port, west port, north port, south port and local port.



Figure 3.3 Router architecture in 2D configuration

Moreover all the ports has bidirectional nature, can be configured as input and output to send and receive the data. The ports are named as: east input port, north input port,west input port, south input port and local input port as inputs, east input port, north input port,west input port, south input port, as outputs. The 2D NoC router (5×5) crossbar switch is shown in the fig. 3.2 which depicts the input and output inputs of the 2D NoC. The architecture of 2D router is shown in fig.3.3. The

data is coming from east input port, west input port, north input port, south input port and local input ports in the packet form and connected to the control unit which decide the size of the data packets need to transmit. The data coming from the east output port, west output port, north output port, south output port and local output port is stored in their associated registers. The design and hardware implementation of the router is done based on the policies, routing and deal with packet collision. The router is also consisting of logic blocks which is helpful in the implementation of the flow control, routing, policies and details the full strategy for the data transfer in the NoC.

3.2 NoC Router with 3D

The crossbar switch (7×7) for 3D NoC is shown in fig. 3.4. The 3D router accepts the seven inputs and seven outs as input and output directional ports.





The data is coming from down input ports, south input port, east input port, north input port, local input port, Up input port, west input port in the packet form and connected to the control unit which decide the size of the data packets need to transmit. The data coming from the east output port, west output port, north output port, south output port, local output port Up input port, Down input ports is stored in their associated registers. The architecture of 3D NoC Router is shown in fig. 3.5 that depicts the functionality of the NoC design.



Figure 3.5 The Architecture of 3D router

3.2.1 Flow Control

A flow control is the method for packet movement along the NoC because it is involved at both level NoC chip level and local router level. It is possible to do flow control and judge the deadlock free routing for specific measures to avoid certain paths within the NoC. The optimization in the NoC depends on the channel requirements and bandwidth and it can guarantees the requirements and need of flow control. The selection technique of the routing algorithm reduces the critical path and traffic congestion is minimized by implementing virtual connections. The quantity of the communication infrastructure and performance is called the Quality of Service (QoS).

The control mechanism is NoC can be classified as centralized control or distributed control. The routing decisions are taken globally and applicable to all the nodes associated in NoC, followed with a strategy that guarantees that traffic contention is not affecting the NoC performance. The approach avoids the bus requirements that all nodes are sharing in bus in a common time. Time Division Multiplexing (TDM) approach is used in which each packet is concern to frame. The NoC uses the distributed approach in which each router can take the decision locally. Virtual channels are very important for the flow control in NoC. These channels are helpful in multiplexing a single physical channel over different or separate logical channels associated in single and independent buffer queues. The actual use of the VC is to implement the NoC to enhance the performance to avoid deadlocks, increase traffic handling capacity, and optimize wire usage. The situation of the deadlock may occur when multiple request sate coming to service node and network resources are fully busy and the nodes are waiting to each other to be that the connection will be free or nodes are waiting to release the connection and proceeding with the communication entity. It is also possible when the two communicating paths are blocked in cyclic manner. If the status of the resources is changing time to time, live lock may occur but no guarantees of communication may be successful.

3.2.2 Routing

Routing in NoC is very important to suggest the optimal path. The routing algorithm, decides the output port for the forward packet when it arrives form the source router to the destination router. The destination port is decided based in the routing information carried by the packet header. There are many routing techniques used to address in NoC, each one having several tradeoff between cost and performance. The same path is used by the data packet in case of deterministicrouting, when communicating between two particular nodes. The general deterministic routing schemes used are XY routing and source routing. Source core are used to specify the route to destination in source routing. In case of the XY routing technique, the row and column approach for routing is used in which the packet is moving towards rows first and then towards column or vice versa is also possible. In another routing technique, alternative paths are possible to communicate among different nodes, in case of original link or path is not available or route is congested. This routing is referred as adaptive routing. This routing is used in modular approach based design for large scale networks in which multiple chances may occur of failure the original links and traffic may be delayed or blocked due to the discontinuation of the original link.

The link load is evaluated using dynamic evaluation technique and follows the strategy based on dynamic load balancing. The other examples related to adaptive routing methods are Negative First (NF) algorithm and West First (WF) algorithm.

The static routing is also one of the important routing that provides the path between the different cores and depends on the time required in the completion of the application, but in case of dynamic routing the routing path is depending on the run time required in the completion of the communication. The data packet can have a single target to deliver the data, called unicast routing whereas one data can be routed to multiple nodes simultaneously with the help of multicast routing technique. It is similar to bus communication. In the same way, a broadcast communication has the destination for all the nodes but narrowcast communication is started by a master node and associated with single slave. The routing methodis also said as minimal routingor non-minimal routing. In case of minimal routing, it is suggested that the shortest path is always decided. In this routing technique a boundary box exist virtually and validates that the decreasing paths for source node to the destination node are valid. On other side, the source to destination distance is increasing based on the non-minimal algorithm. Routing methods and algorithms are very much helpful to prevent the conditions of live locks and deadlocks. The turn model is also one of the routing concept that is also prohibited by certain turns may have increase the risk of deadlock in the network cycle. The network locations are restricted with the help of even odd concept followed in turn model by some types of turns are possible. Another routing algorithm used in networks is hot potato routing. In this technique the data packet is forwarded immediately in the direction of the path which has minimum delay or shortest route for lowest delay. If the packet is not considered by the network, the same packet is returned back to the source router. That's why this routing is also called deflective routing. The packets are not in the buffer and each packet is having multiple set of outputs preferred by the packet itself and used in forward direction against each possibility.

3.2.3 Arbitration

In the NoC communication it is possible that the multiple requests are arriving to the service node. The arbitration logic is used when the routing algorithm decides the output port for packet data transfer and multiple packets are arrived at the source router simultaneously and requesting to provide the service. There are multiple options to implement the logical arbiter. It can becentralized with one per port or distributed one per router. The execution of the operations is taken based on the priority assignment. It may be based on dynamic priority and static priority among the different ports. Router switching matrix is associated with the arbiter and central arbiter can optimize the utilization of router switching matrix. Such arbiter can enhance the latency. Arbitration model is also decided whether the network is following the loss communication model or delayed model. There are the chances of the delaying the packet in the delay model but the packet will not drop out. It is possible that the packet may be dropped in case of loss model. In case of heavy traffic or congestion the packet may be dropped in loss model. In such situations retransmission logic should be implemented. The main function of the crossbar switch is to provide arbitration that provides the solution for conflicting requests occurring for the same output. The speed of the operation directly related to the scheduler and delayed by the arbiter. The designing of the faster arbiter is the challenging and the integration in NoC is very important. The block diagram of the arbiter logic is shown in fig. 3.6. The arbiter has the multiple inputs coming in the form of packets and assigned to the FIFO, logic which assigns the priority. The Requests and grant signals are associated with the crossbar scheduler and directly to arbiter that handles the requests from all the directions in 2D and 3D NoC. Part from that the virtual channel allocation and routing mechanism also the integral part of the arbiter.



3.2.4 Switching

The switching technique decides the way to transmit data from source node to the target node. The entire path is decided already in circuit switching method. The path is included as channel path, routers or whole, established by the header and reserved for sending the complete packet. The payload is not transferred till time the completer path has been reserved. It can enhance the latency. But if the path is decided, the approach can provide some guaranteed throughput. On the other hand, in the packet switching all the flits of data packets are transferred as the header to make the connection among the routers. In this technique the designer has the freedom to decide different forward and buffering strategies than can impact on the overall NoC by implementing the operation concurrently or in parallel pipelined mode by sending the flits. The connection is established to the next routers in pipelined mode. The nodes store the full packet before staring transmission to the next node in defined path. It is the part of store and forward technique in NoC transmission. It should be ensured that the buffer size for each node must be enough to store the full packet

On the other side, in wormhole strategy the nodes have the freedom to take routing decisions and sending the packets to the destination as the header arrives. In the subsequent way, the flits accept the header as it arrives. It reduces the router latency but many links may be blocked due to packet stalling. So there is the involvement of risk factor. Similar to wormhole technique, there is another mechanism called virtual-cut-through in which data is filtered first to the next node in the path before forward and node is waiting for the full packet acceptance by the next node and confirms the complete reception of data. Thus in case of stalling the packets, only current node is affected and other links are safe.

3.2.5 Buffering

The buffering is the methodology used to store the router information in case of traffic congestion and the packet is not reaching directly to the network. The buffer size, buffering strategy and location has very important impact on the traffic intensity. So, the NoC performance is affected by buffering. The router area is also dependent on buffer size. Some routers can have one buffer shared by all the ports and one can have multiple buffers associated with each input and output ports. The biggest advantage of the first concept is that it provides optimized area but the control is very much complex and extra time is required to deal with buffer and overflow. Each port is having its own buffer and FIFO logic is used to implement although other techniques are also available to implement the logic. The techniques may be less efficient because different input ports may require the data storage in single structure.

3.3 NoC Performance Parameters

NoC performance can be evaluated using three parameters.

Bandwidth

Through put Latency

Bandwidth: the bandwidth is called the maximum data rate required to send the messages in network. The unit of the bandwidth is bps. It has the complete packet which includes the bits of header, payload bits and tail bits.

Throughput: throughput is defined as the maximum volume of the incoming traffic considered by the network. It is considered as the maximum amount of information delivered to the load in per unit time. The throughput is measure the message per clock cycle or message per second. One network can have normal throughput which is independent from the size of the message of the network distributed by the size of message and by the length of the network as a result when normal throughput is considered per node or per clock cycle or bit per second. **Latency:** latency is the time required in elapsing between the start of the transmission of the data packet or message and complete receiving at the destination node. Latency is measured with respect to the time unit and mostly following the comparison between different design choices. Latency is also described in terms of simulator clock pulses or clock cycles. In normal way latency is not considered for single packet but average latency is considered to judge the performance of complete network, on the other side when all the messages are presenting larger latency then the average latency become important. That is why the standard deviation is very important to measure against the network.