# 6.8 Introduction to multiprocessor network topologies

(Original section[1])

Multicore chips require on-chip networks to connect cores together, and clusters require local area networks to connect servers together. This section reviews the pros and cons of different interconnection network topologies.

Network costs include the number of switches, the number of links on a switch to connect to the network, the width (number of bits) per link, and length of the links when the network is mapped into silicon. For example, some cores or servers may be adjacent and others may be on the other side of the chip or the other side of the datacenter. Network performance is multifaceted as well. It includes the latency on an unloaded network to send and receive a message, the throughput in terms of the maximum number of messages that can be transmitted in a given time period, delays caused by contention for a portion of the network, and variable performance depending on the pattern of communication. Another obligation of the network may be fault tolerance, since systems may be required to operate in the presence of broken components. Finally, in this era of energy-limited systems, the energy efficiency of different organizations may trump other concerns.

Networks are normally drawn as graphs, with each edge of the graph representing a link of the communication network. In the figures in this section, the processor-memory node is shown as a black square and the switch is shown as a colored circle. We assume here that all links are *bidirectional*; that is, information can flow in either direction. All networks consist of *switches* whose links go to processor-memory nodes and to other switches. The first network connects a sequence of nodes together:



This topology is called a *ring*. Since some nodes are not directly connected, some messages will have to hop along intermediate nodes until they arrive at the final destination.

Unlike a bus—a shared set of wires that allows broadcasting to all connected devices—a ring is capable of many simultaneous transfers.

Because there are numerous topologies to choose from, performance metrics are needed to distinguish these designs. Two are popular. The first is total *network bandwidth*, which is the bandwidth of each link multiplied by the number of links. This represents the peak bandwidth. For the ring network above, with $P$ processors, the total network bandwidth would be $P$ times the bandwidth of one link; the total network bandwidth of a bus is just the bandwidth of that bus.

> **Network bandwidth**: Informally, the peak transfer rate of a network; can refer to the speed of a single link or the collective transfer rate of all links in the network.

To balance this best bandwidth case, we include another metric that is closer to the worst case: the *bisection bandwidth*. This metric is calculated by dividing the machine into two halves. Then you sum the bandwidth of the links that cross that imaginary dividing line. The bisection bandwidth of a ring is two times the link bandwidth. It is one times the link bandwidth for the bus. If a single link is as fast as the bus, the ring is only twice as fast as a bus in the worst case, but it is $P$ times faster in the best case.

> **Bisection bandwidth**: The bandwidth between two equal parts of a multiprocessor. This measure is for a worst case split of the multiprocessor.

Since some network topologies are not symmetric, the question arises of where to draw the imaginary line when bisecting the machine. Bisection bandwidth is a worst-case metric, so the answer is to choose the division that yields the most pessimistic network performance. Stated alternatively, calculate all possible bisection bandwidths and pick the smallest. We take this pessimistic view because parallel programs are often limited by the weakest link in the communication chain.

At the other extreme from a ring is a *fully connected network*, where every processor has a bidirectional link to every other processor. For fully connected networks, the total network bandwidth is $P \times (P - 1)/2$, and the bisection bandwidth is $(P/2)^2$.
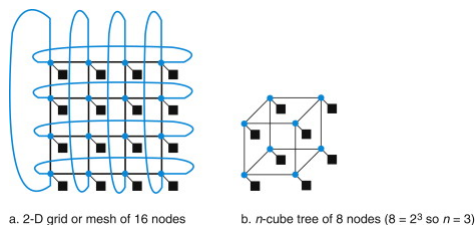
> **Fully connected network**: A network that connects processor-memory nodes by supplying a dedicated communication link between every node.

The tremendous improvement in performance of fully connected networks is offset by the tremendous increase in cost. This consequence inspires engineers to invent new topologies that are between the cost of rings and the performance of fully connected networks. The evaluation of success depends in large part on the nature of the communication in the workload of parallel programs run on the computer.

The number of different topologies that have been discussed in publications would be difficult to count, but only a few have been used in commercial parallel processors. The figure below illustrates two of the popular topologies.

## Figure 6.8.1: Network topologies that have appeared in commercial parallel processors (COD Figure 6.14).

The colored circles represent switches and the black squares represent processor-memory nodes. Even though a switch has many links, generally only one goes to the processor. The Boolean $n$-cube topology is an $n$-dimensional interconnect with $2^n$ nodes, requiring $n$ links per switch (plus one for the processor) and thus $n$ nearest-neighbor nodes. Frequently, these basic topologies have been supplemented with extra arcs to improve performance and reliability.



a. 2-D grid or mesh of 16 nodes          b. $n$-cube tree of 8 nodes ($8 = 2^3$ so $n = 3$)

An alternative to placing a processor at every node in a network is to leave only the switch at some of these nodes. The switches are smaller than processor-memory-switch nodes, and thus may be packed more densely, thereby lessening distance and increasing performance. Such networks are frequently called *multistage networks* to reflect the multiple steps that a message may travel. Types of multistage networks are as numerous as single-stage networks; the figure below illustrates two of the popular multistage organizations. A
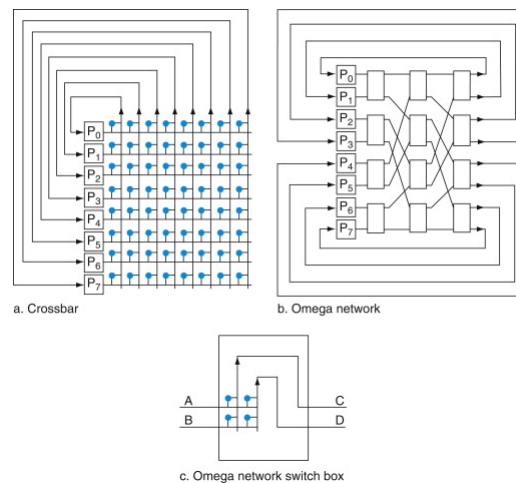
*fully connected* or *crossbar network* allows any node to communicate with any other node in one pass through the network. An Omega network uses less hardware than the crossbar network ($2n \ log_2 n$ versus $n^2$ switches), but contention can occur between messages, depending on the pattern of communication. For example, the Omega network in the figure below cannot send a message from $P_0$ to $P_6$ at the same time that it sends a message from $P_1$ to $P_4$.

**Multistage network**: A network that supplies a small switch at each node.

**Crossbar network**: A network that allows any node to communicate with any other node in one pass through the network.

Figure 6.8.2: Popular multistage network topologies for eight nodes (COD Figure 6.15).

The switches in these drawings are simpler than in earlier drawings because the links are unidirectional; data come in at the left and exit out the right link. The switch box in c can pass A to C and B to D or B to C and A to D. The crossbar uses $n^2$ switches, where $n$ is the number of processors, while the Omega network uses $2n \ log_2 n$ of the large switch boxes, each of which is logically composed of four of the smaller switches. In this case, the crossbar uses 64 switches versus 12 switch boxes, or 48 switches, in the Omega network. The crossbar, however, can support any combination of messages between processors, while the Omega network cannot.



a. Crossbar

b. Omega network

c. Omega network switch box

**Implementing network topologies**

This simple analysis of all the networks in this section ignores important practical considerations in the construction of a network. The distance of each link affects the cost of communicating at a high clock rate—generally, the longer the distance, the more expensive it is to run at a high clock rate. Shorter distances also make it easier to assign more wires to the link, as the power to drive many wires is less if the wires are short. Shorter wires are also cheaper than longer wires. Another practical limitation is that the three-dimensional drawings must be mapped onto chips that are essentially two-dimensional media. The final concern is energy. Energy concerns may force multicore chips to rely on simple grid topologies, for example. The bottom line is that topologies that appear elegant when sketched on the blackboard may be impractical when constructed in silicon or in a datacenter.

Now that we understand the importance of clusters and have seen topologies that we can follow to connect them together, we next look at the hardware and software of the interface of the network to the processor.

**Check yourself**

True or false: For a ring with P nodes, the ratio of the total network bandwidth to the bisection bandwidth is P/2.

**Answer:** True.

(*1) This section is in original form.

⚠ **Provide feedback on this section**