

Performance Analysis of k -ary n -cube Interconnection Networks

WILLIAM J. DALLY, MEMBER, IEEE

Abstract—VLSI communication networks are wire-limited. The cost of a network is not a function of the number of switches required, but rather a function of the wiring density required to construct the network. This paper analyzes communication networks of varying dimension under the assumption of constant wire bisection. Expressions for the latency, average case throughput, and hot-spot throughput of k -ary n -cube networks with constant bisection are derived that agree closely with experimental measurements. It is shown that low-dimensional networks (e.g., tori) have lower latency and higher hot-spot throughput than high-dimensional networks (e.g., binary n -cubes) with the same bisection width.

Index Terms—Communication networks, concurrent computing, interconnection networks, message-passing multiprocessors, parallel processing, VLSI.

I. INTRODUCTION

THE critical component of a concurrent computer is its communication network. Many algorithms are communication rather than processing limited. Fine-grain concurrent programs execute as few as ten instructions in response to a message [7]. To efficiently execute such programs the communication network must have a latency no greater than about ten instruction times, and a throughput sufficient to permit a large fraction of the nodes to transmit simultaneously. Low-latency communication is also critical to support code sharing and garbage collection across nodes.

As the grain size of concurrent computers continues to decrease, communication latency becomes a more important factor. The diameter of the machine grows, messages are sent more frequently, and fewer instructions are executed in response to each message. Low latency is more difficult to achieve in a fine-grain machine because the available wiring space grows more slowly than the expected traffic. Since the machine must be constructed in three dimensions, the bisection area grows only as $N^{2/3}$ while traffic grows at least as fast as N , the number of nodes.

Manuscript received September 3, 1987; revised March 28, 1988. This work was supported in part by the Defense Advanced Research Projects Agency under Contracts N000014-80-C-0622 and N00014-85-K-0124 and in part by a National Science Foundation Presidential Young Investigator Award with matching funds from General Electric Corporation. A preliminary version of this paper appeared in the Proceedings of the 1987 Stanford Conference on Advanced Research in VLSI [10].

The author is with the Artificial Intelligence Laboratory and the Laboratory for Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139.

IEEE Log Number 9034541.

VLSI systems are wire-limited. The cost of these systems is predominantly that of connecting devices, and the performance is limited by the delay of these interconnections. Thus, to achieve the required performance, the network must make efficient use of the available wire. The topology of the network must map into the three physical dimensions so that messages are not required to *double back* on themselves, and in a way that allows messages to use all of the available bandwidth along their path.

This paper considers the problem of constructing *wire-efficient* communication networks, networks that give the optimum performance for a given wire density. We compare networks holding wire bisection, the number of wires crossing a cut that evenly divides the machine, constant. Thus, we compare low-dimensional networks with wide communication channels against high-dimensional networks with narrow channels. We investigate the class of k -ary n -cube interconnection networks and show that low-dimensional networks outperform high-dimensional networks with the same bisection width.

The remainder of this paper describes the design of wire-efficient communication networks. Section II describes the assumptions on which this paper is based. The family of k -ary n -cube networks is described in Section II-A. We restrict our attention to k -ary n -cubes because it is the dimension of the network that is important, not the details of its topology. Section II-B introduces *wormhole routing* [20], a low-latency routing technique. Network cost is determined primarily by wire density which we will measure in terms of bisection width. Section II-C introduces the idea of *bisection width*, and discusses delay models for network channels. A performance model of these networks is derived in Section III. Expressions are given for network latency as a function of traffic that agree closely with experimental results. Under the assumption of constant wire density, it is shown that low-dimensional networks achieve lower latency and better hot-spot throughput than do high-dimensional networks.

II. PRELIMINARIES

A. k -ary n -cubes

Many different network topologies have been proposed for use in concurrent computers: trees [6], [15], [21], Benes networks [4], Batcher sorting networks [2], shuffle exchange networks [23], *Omega* networks [14], *indirect* binary n -cube or *flip* networks [3], [22], and direct binary n -cubes [19], [17], [24]. The binary n -cube is a special case of the family of k -

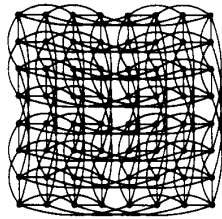


Fig. 1. A binary 6-cube embedded in the plane.

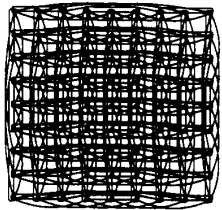


Fig. 2. A ternary 4-cube embedded in the plane.

ary n -cubes, cubes with n dimensions and k nodes in each dimension.

Most concurrent computers have been built using networks that are either k -ary n -cubes or are isomorphic to k -ary n -cubes: rings, meshes, tori, direct and indirect binary n -cubes, and Omega networks. Thus, in this paper we restrict our attention to k -ary n -cube networks. We refer to n as the *dimension* of the cube and k as the *radix*. Dimension, radix, and number of nodes are related by the equation

$$N = k^n, \quad (k = \sqrt[n]{N}, n = \log_k N). \quad (1)$$

It is the dimension of the network that is important, not the details of its topology.

A node in the k -ary n -cube can be identified by n -digit radix k address, a_0, \dots, a_{n-1} . The i th digit of the address, a_i , represents the node's position in the i th dimension. Each node can forward messages to its upper neighbor in each dimension, i , with address $a_0, \dots, a_i + 1 \pmod{k}, \dots, a_{n-1}$.

In this paper, we assume that our k -ary n -cubes are unidirectional for simplicity. We will see that our results do not change appreciably for bidirectional networks. For an actual machine, however, there are many compelling reasons to make our networks bidirectional. Most importantly, bidirectional networks allow us to exploit locality of communication. If an object A sends a message to an object B , there is a high probability of B sending a message back to A . In a bidirectional network, a roundtrip from A to B can be made short by placing A and B close together. In a unidirectional network, a roundtrip will always involve completely circling the machine in at least one dimension.

Figs. 1-3 show three k -ary n -cube networks in order of decreasing dimension. Fig. 1 shows a binary 6-cube (64 nodes). A 3-ary 4-cube (81 nodes) is shown in Fig. 2. An 8-ary 2-cube (64 nodes), or torus, is shown in Fig. 3. Each line in Fig. 1 represents two communication channels, one in each direction, while each line in Figs. 2 and 3 represents a single communication channel.

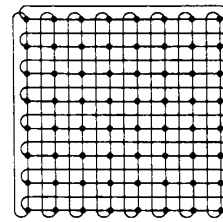


Fig. 3. An 8-ary 2-cube (torus).

B. Wormhole Routing

In this paper, we consider networks that use *wormhole* [20] rather than *store-and-forward* [25] routing. Instead of storing a packet completely in a node and then transmitting it to the next node, wormhole routing operates by advancing the head of a packet directly from incoming to outgoing channels. Only a few flow control digits (flits) are buffered at each node. A *flit* is the smallest unit of information that a queue or channel can accept or refuse.

As soon as a node examines the header flit(s) of a message, it selects the next channel on the route and begins forwarding flits down that channel. As flits are forwarded, the message becomes spread out across the channels between the source and destination. It is possible for the first flit of a message to arrive at the destination node before the last flit of the message has left the source. Because most flits contain no routing information, the flits in a message must remain in contiguous channels of the network and cannot be interleaved with the flits of other messages. When the header flit of a message is blocked, all of the flits of a message stop advancing and block the progress of any other message requiring the channels they occupy.

A method similar to wormhole routing, called *virtual cut-through*, is described in [13]. Virtual cut-through differs from wormhole routing in that it buffers messages when they block, removing them from the network. With wormhole routing, blocked messages remain in the network.

Fig. 4 illustrates the advantage of wormhole routing. There are two components of latency, distance and message aspect ratio. The distance D is the number of *hops* required to get from the source to the destination. The message aspect ratio (message length L normalized to the channel width W) is the number of channel cycles required to transmit the message across one channel. The top half of the figure shows store-and-forward routing. The message is entirely transmitted from node N_0 to node N_1 , then from N_1 to N_2 and so on. With store-and-forward routing, latency is the product of D and L/W .

$$T_{SF} = T_c \left(D \times \frac{L}{W} \right). \quad (2)$$

The bottom half of Fig. 4 shows wormhole routing. As soon as a flit arrives at a node, it is forwarded to the next node. With wormhole routing, latency is reduced to the sum of D and L/W .

$$T_{WH} = T_c \left(D + \frac{L}{W} \right). \quad (3)$$

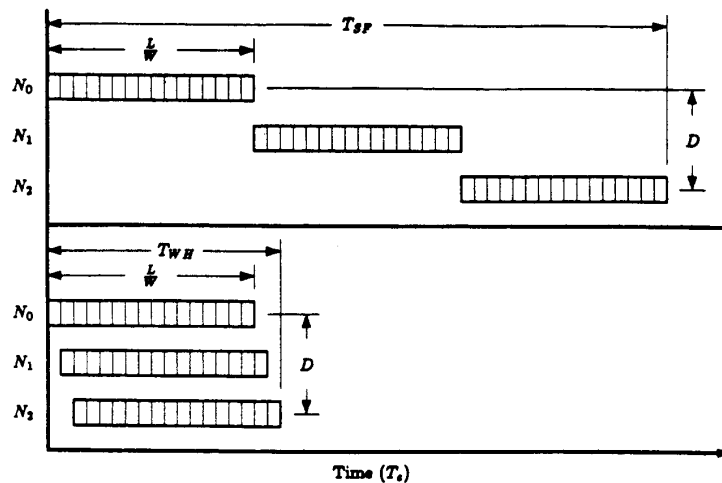


Fig. 4. Latency of store-and-forward routing (top) versus wormhole routing (bottom).

In both of these equations, T_c is the channel cycle time, the amount of time required to perform a transaction on a channel.

C. VLSI Complexity

VLSI computing systems [16] are wire-limited; the complexity of what can be constructed is limited by wire density, the speed at which a machine can run is limited by wire delay, and the majority of power consumed by a machine is used to drive wires. Thus, machines must be organized both logically and physically to keep wires short by exploiting locality wherever possible. The VLSI architect must organize a computing system so that its form (physical organization) fits its function (logical organization).

Networks have traditionally been analyzed under the assumption of constant channel bandwidth. Under this assumption each channel is one bit wide ($W = 1$) and has unit delay ($T_c = 1$). The constant bandwidth assumption favors networks with high dimensionality (e.g., binary n -cubes) over low-dimensional networks (e.g., tori). This assumption, however, is not consistent with the properties of VLSI technology. Networks with many dimensions require more and longer wires than do low-dimensional networks. Thus, high-dimensional networks cost more and run more slowly than low-dimensional networks. A realistic comparison of network topology must take both wire density and wire length into account.

To account for wire density, we will use bisection width [26] as a measure of network cost. The bisection width of a network is the minimum number of wires cut when the network is divided into two equal halves. Rather than comparing networks with constant channel width W , we will compare networks with constant bisection width. Thus, we will compare low-dimensional networks with large W with high-dimensional networks with small W .

The delay of a wire depends on its length l . For short wires, the delay t_s is limited by charging the capacitance of the wire and varies logarithmically with wire length.

$$t_s = \tau_{inv} e \log_e Kl \quad (4)$$

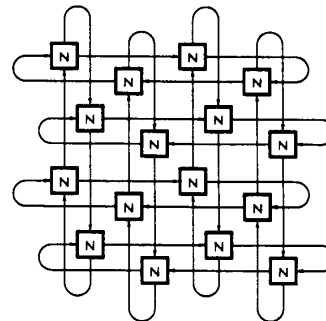


Fig. 5. A folded torus system.

where τ_{inv} is the inverter delay, and K is a constant depending on capacitance ratios. For long wires, delay t_l is limited by the speed of light.

$$t_l = \frac{l\sqrt{\epsilon_r}}{c} \quad (5)$$

In this paper, we will consider three delay models: constant delay T_c independent of length, logarithmic delay $T_c \propto \log l$, and linear delay $T_c \propto l$. Our main result, that latency is minimized by low-dimensional networks, is supported by all three models.

III. PERFORMANCE ANALYSIS

In this section, we compare the performance of unidirectional k -ary n -cube interconnection networks using the following assumptions.

- Networks must be embedded into the plane. If a three-dimensional packaging technology becomes available, the comparison changes only slightly.

- Nodes are placed systematically by embedding $n/2$ logical dimensions in each of the two physical dimensions. We assume that both n and k are even integers. The long end-around connections shown in Fig. 3 can be avoided by folding the network as shown in Fig. 5.

- For networks with the same number of nodes, *wire density is held constant*. Each network is constructed with the same bisection width B , the total number of wires crossing the midpoint of the network. To keep the bisection width constant, we vary the width W of the communication channels. We normalize to the bisection width of a bit-serial ($W = 1$) binary n -cube.

- The networks use *wormhole* routing.

- Channel delay T_c is a function of wire length l . We begin by considering channel delay to be constant. Later, the comparison is performed for both logarithmic and linear wire delays; $T_c \propto \log l$ and $T_c \propto l$.

When k is even, the channels crossing the midpoint of the network are all in the highest dimension. For each of the \sqrt{N} rows of the network, there are $k^{((n/2)-1)}$ of these channels in each direction for a total of $2\sqrt{N}k^{((n/2)-1)}$ channels. Thus, the bisection width B of a k -ary n -cube with W -bit wide communication channels is

$$B(k, n) = 2W\sqrt{N}k^{((n/2)-1)} = \frac{2WN}{k}. \quad (6)$$

For a binary n -cube, $k = 2$, the bisection width is $B(2, n) = WN$. We set B equal to N to normalize to a binary n -cube with unit width channels, $W = 1$. The channel width $W(k, n)$ of a k -ary n -cube with the same bisection width B follows from (6):

$$\frac{2W(k, n)N}{k} = N,$$

$$W(k, n) = \frac{k}{2}. \quad (7)$$

The peak wire density is greater than the bisection width in networks with $n > 2$ because the lower dimensions contribute to wire density. The maximum density, however, is bounded by

$$\begin{aligned} D_{\max} &= 2W\sqrt{N} \sum_{i=0}^{\frac{n}{2}-1} k^i = k\sqrt{N} \sum_{i=0}^{\frac{n}{2}-1} k^i \\ &= k\sqrt{N} \left(\frac{k^{n/2} - 1}{k - 1} \right) \\ &= k\sqrt{N} \left(\frac{\sqrt{N} - 1}{k - 1} \right) < \left(\frac{k}{k - 1} \right) B. \end{aligned} \quad (8)$$

A plot of wire density as a function of position for one row of a binary 20-cube is shown in Fig. 6. The density is very low at the edges of the cube and quite dense near the center. The peak density for the row is 1364 at position 341. Compare this density to the bisection width of the row, which is 1024. In contrast, a two-dimensional torus has a wire density of 1024 independent of position. One advantage of high-radix networks is that they have a very uniform wire density. They make full use of available area.

Each processing node connects to $2n$ channels (n input and n output) each of which is $k/2$ bits wide. Thus, the number

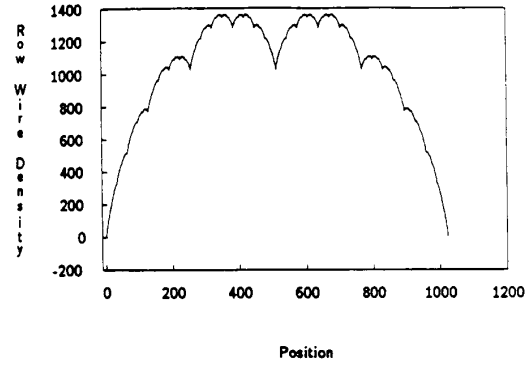


Fig. 6. Wire density versus position for one row of a binary 20-cube.

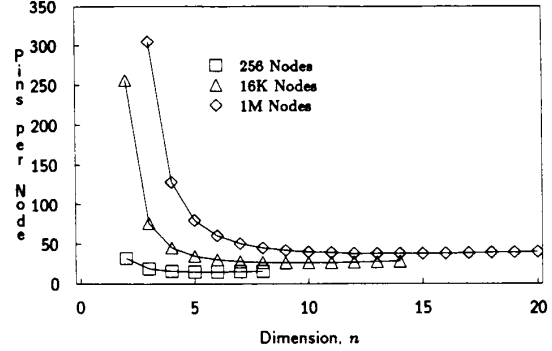


Fig. 7. Pin density versus dimension for 256, 16K, and 1M nodes.

of pins per processing node is

$$N_p = nk. \quad (9)$$

A plot of pin density as a function of dimension for $N = 256$, 16K, and 1M nodes¹ is shown in Fig. 7. Low-dimensional networks have the disadvantage of requiring many pins per processing node. A two-dimensional network with 1M nodes (not shown) requires 2048 pins and is clearly unrealizable. However, the number of pins decreases very rapidly as the dimension n increases. Even for 1M nodes, a dimension 4 node has only 128 pins. All of the configurations that give low latency also give a reasonable pin count.

A. Latency

Latency T_l is the sum of the latency due to the network and the latency due to the processing node,

$$T_l = T_{\text{net}} + T_{\text{node}}. \quad (10)$$

In this paper, we are concerned only with T_{net} . Techniques to reduce T_{node} are described in [7] and [11].

If we select two processing nodes, P_i, P_j , at random, the average number of channels that must be traversed to send a message from P_i to P_j is given by

$$D = \left(\frac{k-1}{2} \right) n. \quad (11)$$

¹ 1K=1024 and, 1M=1K×1K=1048576.

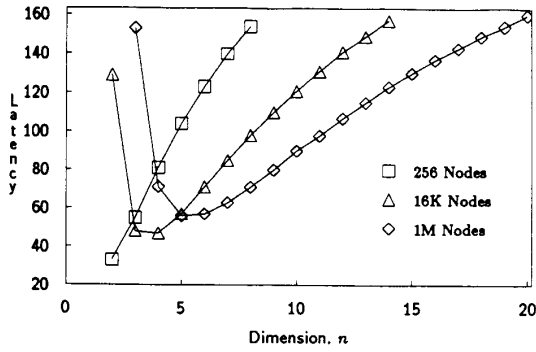


Fig. 8. Latency versus dimension for 256, 16K, and 1M nodes, constant delay.

The average latency of a k -ary n -cube is calculated by substituting (7) and (11) into (3)

$$T_{\text{net}} = T_c \left(\left(\frac{k-1}{2} \right) n + \frac{2L}{k} \right). \quad (12)$$

Fig. 8 shows the average network latency T_{net} as a function of dimension n for k -ary n -cubes with 2^8 (256), 2^{14} (16K), and 2^{30} (1M) nodes.² The left most data point in this figure corresponds to a torus ($n = 2$) and the rightmost data point corresponds to a binary n -cube ($k = 2$). This figure assumes constant wire delay T_c and a message length L of 150 bits. This choice of message length was based on the analysis of a number of fine-grain concurrent programs [7]. Although constant wire delay is unrealistic, this figure illustrates that even ignoring the dependence of wire delay on wire length, low-dimensional networks achieve lower latency than high-dimensional networks.

The latency of the tori on the left side of Fig. 8 is limited almost entirely by distance. The latency of the binary n -cubes on the right side of the graph is limited almost entirely by aspect ratio. With bit serial channels, these cubes take 150 cycles to transmit their messages across a single channel.

In an application that exploits locality of communication, the distance between communicating objects is reduced. In such a situation, the latency of the low-dimensional networks, dominated by distance (the left side of Fig. 8) is reduced. High-dimensional networks, on the other hand, cannot take advantage of locality. Their latency, because it is dominated by message length, will remain high.

In applications that send short messages, the component of latency due to message length is reduced resulting in lower latency for high-dimensional networks (the right side of Fig. 8).

For the three cases shown in Fig. 8, minimum latencies are achieved for $n = 2, 4,$ and $5,$ respectively. In general, the lowest latency is achieved when the component of latency due to distance D and the component due to message length L/W

² For the sake of comparison, we allow radix to take on noninteger values. For some of the dimensions considered, there is no integer radix k that gives the correct number of nodes. In fact, this limitation can be overcome by constructing a *mixed-radix cube* [5].

are approximately equal, $D \approx L/W$. The following assertion makes this statement more precise.

Assertion: Minimum latency T_{net} occurs at a dimension $n \leq n_x$, where n_x is the dimension for which $D = L/W$.

Proof: Differentiating (12) with respect to n gives

$$\frac{\partial T_{\text{net}}}{\partial n} = \frac{k-1-k \log k}{2} + \frac{2L \log^2 k}{k \log N}. \quad (13)$$

For $n = n_x$, substituting $D = L/W$ into (7) and (11) gives

$$4L = nk(k-1) = \frac{k(k-1) \log N}{\log k}. \quad (14)$$

Substituting into the derivative (13) gives

$$\begin{aligned} \frac{\partial T_{\text{net}}}{\partial n} \Big|_{n=n_x} &= \frac{k-1-k \log k}{2} + \frac{(k-1) \log k}{2} \\ &= \frac{k-1-\log k}{2}. \end{aligned} \quad (15)$$

For all $k \geq 2$, $(\partial T_{\text{net}}/\partial n)|_{n=n_x} \geq 0$. The derivative is monotonically increasing for $n \leq n_x$. Thus, the minimum latency ($\partial T_{\text{net}}/\partial n = 0$) occurs for $n < n_x$. \square

Empirically, for all networks with $N < 2^{20}$ and integral valued k and n the minimum latency occurs when k and n are chosen so that $|D - (L/W)|$ is minimized.

The longest wire in the system becomes a bottleneck that determines the rate at which each channel operates, T_c . The length of this wire is given by

$$l = k^{(n/2)-1}. \quad (16)$$

If the wires are sufficiently short, delay depends logarithmically on wire length. If the channels are longer, they become limited by the speed of light, and delay depends linearly on channel length. Substituting (16) into (4) and (5) gives

$$T_c \propto \begin{cases} 1 + \log_e l = 1 + \left(\frac{n}{2} - 1 \right) \log_e k & \text{logarithmic delay} \\ l = k^{(n/2)-1} & \text{linear delay.} \end{cases} \quad (17)$$

We substitute (17) into (12) to get the network latency for these two cases:

$$T_l \propto \begin{cases} \left(1 + \left(\frac{n}{2} - 1 \right) \log_e k \right) \left(\left(\frac{k-1}{2} \right) n + \frac{2L}{k} \right) & \text{logarithmic delay} \\ k^{(n/2)-1} \left(\left(\frac{k-1}{2} \right) n + \frac{2L}{k} \right) & \text{linear delay.} \end{cases} \quad (18)$$

Fig. 9 shows the average network latency as a function of dimension for k -ary n -cubes with 2^8 (256), 2^{14} (16K), and 2^{20} (1M) nodes, assuming logarithmic wire delay and a message length, L , of 150. Fig. 10 shows the same data assuming linear wire delays. In both figures, the leftmost data point corresponds to a torus ($n = 2$) and the rightmost data point corresponds to a binary n -cube ($k = 2$).

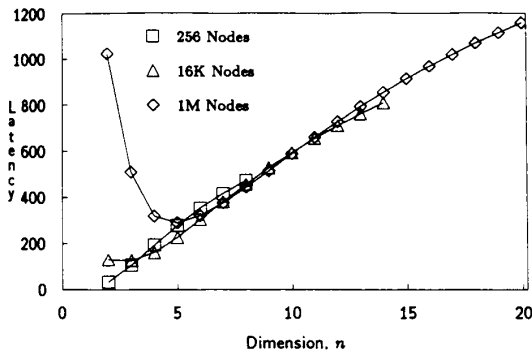


Fig. 9. Latency versus dimension for 256, 16K, and 1M nodes, logarithmic delay.

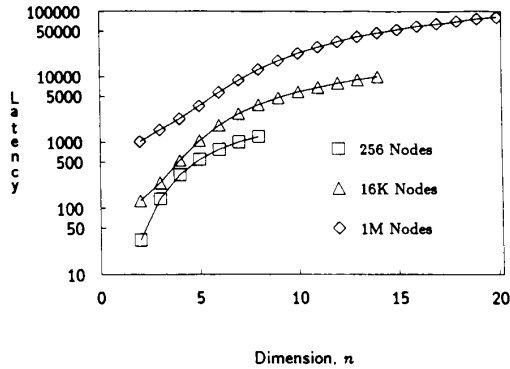


Fig. 10. Latency versus dimension for 256, 16K, and 1M nodes, linear delay.

In the linear delay case, Fig. 10, a torus ($n = 2$) always gives the lowest latency. This is because a torus offers the highest bandwidth channels and the most direct physical route between two processing nodes. Under the linear delay assumption, latency is determined solely by bandwidth and by the physical distance traversed. There is no advantage in having long channels.

Under the logarithmic delay assumption, Fig. 9, a torus has the lowest latency for small networks ($N = 256$). For the larger networks, the lowest latency is achieved with slightly higher dimensions. With $N = 16K$, the lowest latency occurs when n is three.³ With $N = 1M$, the lowest latency is achieved when n is 5. It is interesting that assuming constant wire delay does not change this result much. Recall that under the (unrealistic) constant wire delay assumption, Fig. 8, the minimum latencies are achieved with dimensions of 2, 4, and 5, respectively.

The results shown in Figs. 8–10 were derived by comparing networks under the assumption of constant wire cost to a binary n -cube with $W = 1$. For small networks it is possible to construct binary n -cubes with wider channels, and for large networks (e.g., 1M nodes) it may not be possible to construct a binary n -cube at all. The available wiring area grows as $N^{2/3}$ while the bisection width of a binary n -cube grows as N . In the case of small networks, the comparison

³ In an actual machine, the dimension n would be restricted to be an even integer.

against binary n -cubes with wide channels can be performed by expressing message length in terms of the binary n -cube's channel width, in effect decreasing the message length for purposes of comparison. The net result is the same: lower dimensional networks give lower latency. Even if we perform the 256 node comparison against a binary n cube with $W=16$, the torus gives the lowest latency under the logarithmic delay model, and a dimension 3 network gives minimum latency under the constant delay model. For large networks, the available wire is less than assumed, so the effective message length should be increased, making low-dimensional networks look even more favorable.

In this comparison, we have assumed that only a single bit of information is in transit on each wire of the network at a given time. Under this assumption, the delay between nodes T_c is equal to the period of each node T_p . In a network with long wires, however, it is possible to have several bits in transit at once. In this case, the channel delay T_c is a function of wire length, while the channel period $T_p < T_c$ remains constant. Similarly, in a network with very short wires we may allow a bit to ripple through several channels before sending the next bit. In this case, $T_p > T_c$. Separating the coefficients T_c and T_p , (3) becomes

$$T_{\text{net}} = \left(T_c D + T_p \frac{L}{W} \right). \quad (19)$$

The net effect of allowing $T_c \neq T_p$ is the same as changing the length L by a factor of $\frac{T_p}{T_c}$ and does not change our results significantly.

When wire cost is considered, low-dimensional networks (e.g., tori) offer lower latency than high-dimensional networks (e.g., binary n -cubes). Tori outperform binary n -cubes because they better match form to function. The logical and physical graphs of the torus are identical; thus, messages always travel the minimum distance from source to destination. In a binary n -cube, on the other hand, the fit between form and function is not as good. A message in a binary n -cube embedded into the plane may have to traverse considerably more than the minimum distance between its source and destination.

B. Throughput

Throughput, another important metric of network performance, is defined as the total number of messages the network can handle per unit time. One method of estimating throughput is to calculate the capacity of a network, the total number of messages that can be in the network at once. Typically the maximum throughput of a network is some fraction of its capacity. The network capacity per node is the total bandwidth out of each node divided by the average number of channels traversed by each message. For k -ary n -cubes, the bandwidth out of each node is nW , and the average number of channels traversed is given by (11), so the network capacity per node Γ is given by

$$\Gamma \propto \frac{nW}{D} \propto \frac{n \left(\frac{k}{2} \right)}{\left(\frac{k-1}{2} \right) n} \approx 1. \quad (20)$$

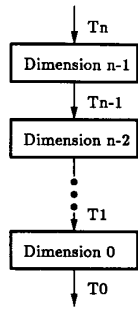


Fig. 11. Contention model for a network.

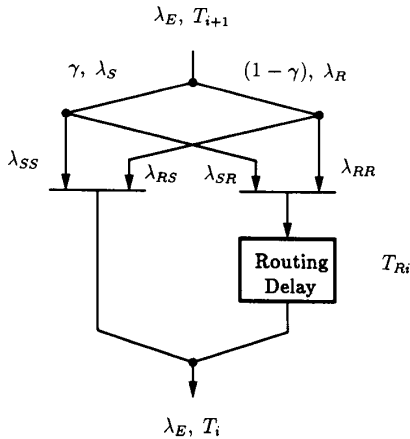


Fig. 12. Contention model for a single dimension.

The network capacity is independent of dimension. For a constant wire density, there is a constant network capacity.

Throughput will be less than capacity because contention causes some channels to block. This contention also increases network latency. To simplify the analysis of this contention, we make the following assumptions.

- Messages are routed using *e*-cube routing (in order of decreasing dimension) [8]. That is, a message at node a_0, \dots, a_{n-1} destined for nodes b_0, \dots, b_{n-1} is first routed in dimension $n - 1$ until it reaches node $a_0, \dots, a_{n-2}, b_{n-1}$. The message is then routed in dimension $n - 2$ until it reaches node $a_0, \dots, a_{n-3}, b_{n-2}, b_{n-1}$, and so on. As shown in Fig. 11, this assumption allows us to consider the contention in each dimension separately.

- The traffic from each node is generated by a Poisson process with arrival rate λ (bits/cycle).

- Message destinations are uniformly distributed and independent.

The arrival rate of λ (bits/cycle) corresponds to $\lambda_E = (\lambda/L)$ (messages/cycles). At the destination, each flit is serviced as soon as it arrives, so the service time at the sink is $T_0 = L/W = 2L/k$. Starting with T_0 we will calculate the service time seen entering each preceding dimension.

For convenience, we will define the following quantities as illustrated in Fig. 12:

$$\gamma = \frac{1}{k} \quad \text{probability that a message skips (does not route in) a dimension,}$$

$$\lambda_S = \gamma\lambda_E \quad \text{rate of traffic that skips the previous dimension, } i+1, \left(\frac{\text{messages}}{\text{cycle}}\right),$$

$$\lambda_R = (1 - \gamma)\lambda_E \quad \text{rate of traffic that routes in the previous dimension, } i + 1, \left(\frac{\text{messages}}{\text{cycle}}\right),$$

$$\lambda_{SS} = \gamma^2\lambda_E \quad \text{rate of traffic that skips both the previous dimension, } i + 1, \text{ and the current dimension } i, \left(\frac{\text{messages}}{\text{cycle}}\right),$$

$$\lambda_{SR} = \gamma(1 - \gamma)\lambda_E \quad \text{rate of traffic that skips the previous dimension, } i + 1, \text{ and routes in the current dimension, } i, \left(\frac{\text{messages}}{\text{cycle}}\right),$$

$$\lambda_{RS} = \gamma(1 - \gamma)\lambda_E \quad \text{rate of traffic that routes in the previous dimension, } i + 1, \text{ and skips the current dimension } i, \left(\frac{\text{messages}}{\text{cycle}}\right),$$

$$\lambda_{RR} = (1 - \gamma)^2\lambda_E \quad \text{rate of traffic that routes in both the previous dimension } i + 1, \text{ and the current dimension } i, \left(\frac{\text{messages}}{\text{cycle}}\right).$$

(21)

Consider a single dimension i of the network as shown in Fig. 12. All messages incur a latency due to contention on entering the dimension. Those messages that are routed incur an additional latency T_{Ri} due to contention during routing. The rate λ_E message stream entering the dimension is composed of two components: a rate λ_S stream that skipped the previous ($i + 1$ st) dimension, and a rate λ_R stream that was routed in the previous dimension. These two streams are in turn split into components that will skip the i th dimension (λ_{SS} and λ_{RS}) and components that will be routed in the i th dimension (λ_{SR} and λ_{RR}). The entering latency seen by one component (say λ_{SS}) is given by multiplying the probability of a collision (in this case $\lambda_{RS}T_i$) by the expected latency due to a collision [in this case $(T_i/2)$]. The components that require routing must also add the latency due to contention during routing, T_{Ri} . Adding up the four components with appropriate weights gives the following equation for T_{i+1} .

$$T_{i+1} = T_i + (1 - \gamma)T_{Ri} + \gamma(1 - \gamma)^3\lambda_E(T_i + T_{Ri})^2 + \gamma^3(1 - \gamma)\lambda_E T_i^2. \quad (22)$$

The first term of (22) is the latency seen entering the next dimension. The second term accounts for the routing latency T_{Ri} incurred by messages routing in this dimension (λ_{SR} and λ_{RR}). The entering latency due to contention when the two routing streams merge is given by the third term. The final term gives the entering latency for the messages that skip the dimension.

For large k , γ is small and the latency is approximated by $T_{i+1} \approx T_i + T_{Ri}$. For $k = 2$ (binary n -cubes), $T_{Ri} = 0$; thus, $T_{i+1} = T_i + (\lambda_E T_i/8)$.

To calculate the routing latency T_{Ri} we use the model shown in Fig. 13. Messages enter the dimension with rate λ_R , route

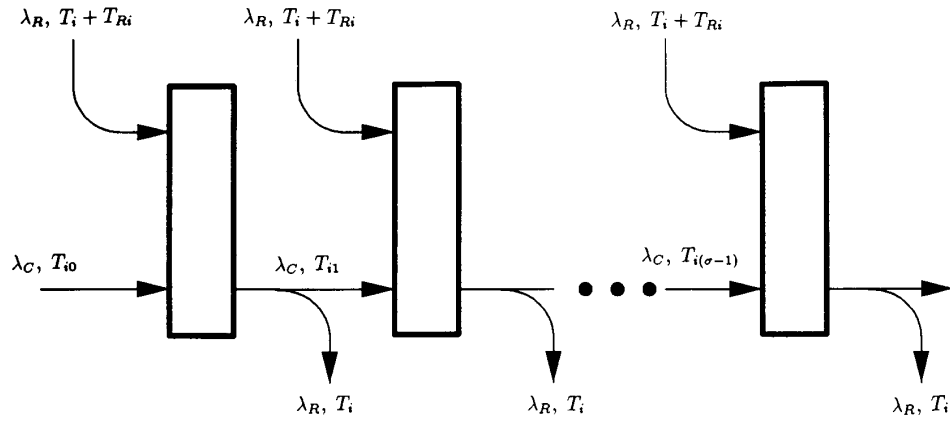


Fig. 13. Contention model for routing latency.

through a number of stages, denoted by boxes, and exit the dimension. The latency due to contention in the stages sums to T_{Ri} . Given that a message is to be routed in a dimension, the expected number of channels traversed by the message is $(k/2)$, one entering channel and $\sigma = (k - 2)/2$ continuing channels. Thus, the *average* message rate on channels continuing in the dimension is $\lambda_C = \sigma\lambda_R$.

Using the average message rate to calculate latency is an approximation. The symmetry of the network assures that the traffic on physical channels is uniform. However, using virtual channels and *e*-cube routing [8] results in logical channels that form a spiral. Traffic on the j th channel on this spiral is given by $\lambda_{Cj} = (j - (i^2 + i)/2k)\lambda_R$. Using the uniform physical message rate results in a slightly pessimistic estimate of latency since contention for the physical channel occurs on flit boundaries while contention for the logical channel occurs on message boundaries.

To compute T_{Ri} we work backwards from the output. The service time in the last continuing channel in dimension i is $T_{i(\sigma-1)} = T_i$. Once we know the service time for the j th channel, T_{ij} , the additional service time due to contention at the $j - 1$ st channel is given by multiplying the probability of a collision $\lambda_R T_{i0}$ by the expected waiting time for a collision $T_{i0}/2$. Repeating this calculation σ times gives us T_{i0} .

$$\begin{aligned} T_{i(j-1)} &= T_{ij} + \frac{\lambda_R T_{i0}^2}{2}, \\ T_{i0} &= T_i + \frac{\sigma\lambda_R T_{i0}^2}{2} = T_i + \frac{\lambda_C T_{i0}^2}{2}, \\ &= \frac{1 - \sqrt{1 - 2\lambda_C T_i}}{\lambda_C}. \end{aligned} \quad (23)$$

Equation (23) is valid only when $\lambda_C < T_i/2$. If the message rate is higher than this limit, there is no steady-state solution and latency becomes infinite. There are two solutions to (23). Here we consider only the smaller of the two latencies. The larger solution corresponds to a state that is not encountered during normal operation of a network.

To calculate T_{Ri} we also need to consider the possibility of a collision on the entering channel.

$$T_{Ri} = T_{i0} \left(1 + \frac{\lambda_C T_{i0}}{2} \right) - T_i. \quad (24)$$

If sufficient queueing is added to each network node, the service times do not increase, only the latency and (24) and (22) become

$$T_{Ri} = \left(\frac{T_i}{1 - \frac{\lambda_C T_0}{2}} \right) \left(1 + \frac{\lambda_C T_0}{2} \right) - T_i, \quad (25)$$

$$T_{i+1} = T_i + (1 - \gamma)T_{Ri} + (\gamma(1 - \gamma)^3 + \gamma^3(1 - \gamma))\lambda_E T_0. \quad (26)$$

To be effective, the total queueing between the source and destination should be greater than the expected increase in latency due to blocking. Two flits of queueing per stage are sufficient when $\lambda < 0.3$ and $L < 200$. Longer messages result in a longer service time T_0 and require additional queueing. The analysis here is pessimistic in that it assumes no queueing. Using (22), we can determine 1) the maximum throughput of the network and 2) how network latency increases with traffic.

Figs. 14 and 15 show how latency increases as a function of applied traffic for 1K node and 4K node k -ary n -cubes. The vertical axis shows latency in cycles. The horizontal axis is traffic per node, λ , in bits/cycle. The figures compare measurements from a network simulator (points) to the latency predicted by (24) (lines). The simulation agrees with the prediction within a few percent until the network approaches saturation.

For 1K networks, a 32-ary 2-cube always gives the lowest latency. For 4K networks, a 16-ary 3-cube gives the lowest latency when $\lambda < 0.2$. Because latency increases more slowly for two-dimensional networks, a 64-ary 2-cube gives the lowest latency when $\lambda > 0.2$. At the left side of each graph ($\lambda = 0$), latency is given by (12). As traffic is applied to the

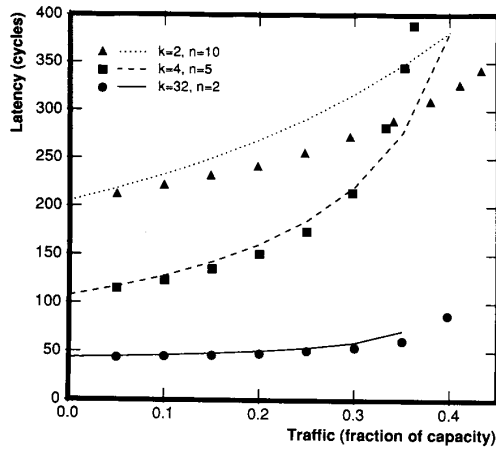


Fig. 14. Latency versus traffic (λ) for 1K node networks: 32-ary 2-cube, 4-ary 5-cube, and binary 10-cube, $L = 200$ bits. Solid line is predicted latency, points are measurements taken from a simulator.

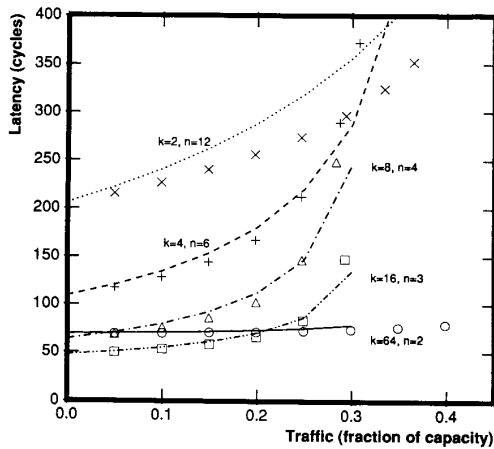


Fig. 15. Latency versus traffic (λ) for 4K node networks: 64-ary 2-cube, 16-ary 3-cube, 8-ary 4-cube, 4-ary 6-cube, and binary 12-cube, $L = 200$ bits. Solid line is predicted latency, points are measurements taken from a simulator.

network, latency increases slowly due to contention in the network until saturation is reached. Saturation occurs when λ is between 0.3 and 0.5 depending on the network topology. Networks should be designed to operate on the flat portion of the curve ($\lambda < 0.25$).

When the network saturates, throughput levels off as shown in Figs. 16 and 17. These figures show how much traffic is delivered (vertical axis) when the nodes attempt to inject a given amount of traffic (horizontal axis). The curve is linear (actual = attempted) until saturation is reached. From this point on, actual traffic is constant. This plateau occurs because 1) the network is source queued, and 2) messages that encounter contention are blocked rather than aborted. In networks where contention is resolved by dropping messages, throughput usually decreases beyond saturation.

To find the maximum throughput of the network, the source service time T_0 is set equal to the reciprocal of the message rate, λ_E , and (22), (23), and (24) are solved for λ_E . At this

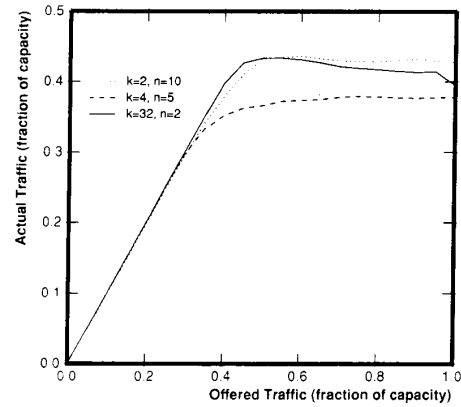


Fig. 16. Actual traffic versus attempted traffic for 1K node networks: 32-ary 2-cube, 4-ary 5-cube, and binary 10-cube, $L = 200$ bits.

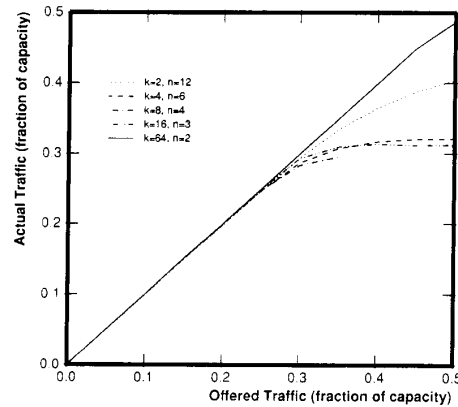


Fig. 17. Actual traffic versus attempted traffic for 4K node networks: 64-ary 2-cube, 16-ary 3-cube, 8-ary 4-cube, 4-ary 6-cube, and binary 12-cube, $L = 200$ bits.

TABLE I
MAXIMUM THROUGHPUT AS A FRACTION OF CAPACITY AND BLOCKING LATENCY IN CYCLES

| Parameter | 1K Nodes | | | 4K Nodes | | | | |
|-------------------------|----------|------|------|----------|------|------|------|------|
| Dimension | 2 | 5 | 10 | 2 | 3 | 4 | 6 | 12 |
| radix | 32 | 4 | 2 | 64 | 16 | 8 | 4 | 2 |
| Max Throughput | 0.36 | 0.41 | 0.43 | 0.35 | 0.31 | 0.31 | 0.36 | 0.41 |
| Latency $\lambda = 0.1$ | 46.1 | 128. | 233. | 70.7 | 55.2 | 79.9 | 135. | 241. |
| Latency $\lambda = 0.2$ | 50.5 | 161. | 269. | 73.1 | 70.3 | 112. | 181. | 288. |
| Latency $\lambda = 0.3$ | 59.3 | 221. | 317. | 78.6 | 135. | 245. | 287. | 357. |

operating point the network can accept no more traffic. Messages are being offered as fast as the network can deliver them. The maximum throughput as a fraction of capacity for k -ary n -cubes with 1K and 4K nodes is tabulated in Table I. Also shown is the total latency for $L = 200$ bit messages at several message rates. The table shows that the additional latency due to blocking is significantly reduced as dimension is decreased.

In networks of constant bisection width, the latency of low-dimensional networks increases more slowly with applied traffic than the latency of high-dimensional networks. At $\lambda = 0.2$,

the 32-ary 2-cube has $\approx \frac{1}{5}$ the latency of the binary 10-cube. At this point, the additional latency due to contention in the 32-ary 2-cube is $7T_c$ compared to $64T_c$ in the binary 10-cube.

Low-dimensional networks handle contention better because they use fewer channels of higher bandwidth and thus get better queueing performance. The shorter service times, L/W , of these networks results in both a lower probability of collision, and a lower expected waiting time in the event of a collision. Thus, the blocking latency at each node is reduced quadratically as k is increased. Low-dimensional networks require more hops, $D = (n(k-1)/2)$, and have a higher rate on continuing channels, λ_c . However, messages travel on the continuing channels more frequently than on the entering channels, thus most contention is with the lower rate channels. Having fewer channels of higher bandwidth also improves hot-spot throughput as described below.

C. Hot-Spot Throughput

In many situations traffic is not uniform, but rather is concentrated into *hot spots*. A *hot spot* is a pair of nodes that accounts for a disproportionately large portion of the total network traffic. As described by Pfister [18] for a shared-memory computer, hot-spot traffic can degrade performance of the entire network by causing congestion.

The *hot-spot throughput* of a network is the maximum rate at which messages can be sent from one specific node P_i to another specific node P_j . For a k -ary n -cube with deterministic routing, the hot-spot throughput, Θ_{HS} , is just the bandwidth of a single channel W . Thus, under the assumption of constant wire cost we have

$$\Theta_{HS} = W = k - 1. \quad (27)$$

Low-dimensional networks have greater channel bandwidth and thus have greater hot-spot throughput than do high-dimensional networks. Low-dimensional networks operate better under nonuniform loads because they do more resource sharing. In an interconnection network the resources are wires. In a high-dimensional network, wires are assigned to particular dimensions and cannot be shared between dimensions. For example, in a binary n -cube it is possible for a wire to be saturated while a physically adjacent wire assigned to a different dimension remains idle. In a torus all physically adjacent wires are combined into a single channel that is shared by all messages that must traverse the physical distance spanned by the channel.

IV. CONCLUSION

Under the assumption of constant wire bisection, low-dimensional networks with wide channels provide lower latency, less contention, and higher hot-spot throughput than higher-dimensional networks with narrow channels. Minimum network latency is achieved when the network radix k and dimension n are chosen to make the components of latency due to distance D and aspect ratio L/W approximately equal. The minimum latency occurs at a very low dimension, 2 for up to 1024 nodes.

Low-dimensional networks reduce contention because having a few high-bandwidth channels results in more resource

sharing and thus better queueing performance than having many low-bandwidth channels. While network capacity and worst-case blocking latency are independent of dimension, low-dimensional networks have a higher maximum throughput and lower average blocking latency than do high-dimensional networks. Improved resource sharing also gives low-dimensional networks higher hot-spot throughput than high-dimensional networks.

The results of this paper have all been made under the assumption of constant channel delay, independent of channel length. The main result, that low-dimensional networks give minimum latency, however, does not change appreciably when logarithmic or linear delay models are considered. In choosing a delay model one must consider how the delay of a switching node compares to the delay of a wire. Current VLSI routing chips [9] have delays of tens of nanoseconds, enough time to drive several meters of wire. For such systems a constant delay model is adequate. As chips get faster and systems get larger, however, a linear delay model will more accurately reflect system performance.

Fat-tree networks have been shown to be universal in the sense that they can *efficiently* simulate any other network of the same volume [15]. However, the analysis of these networks has not considered latency. k -ary n -cubes with appropriately chosen radix and dimension are also universal in this sense. A detailed proof is beyond the scope of this paper. Intuitively, one cannot do any better than to fill each of the three physical dimensions with wires and place switches at every point of intersection. Any point-to-point network can be embedded into such a 3-D mesh with no more than a constant increase in wiring length.

This paper has considered only *direct* networks [19]. The results do not apply to *indirect* networks. The depth and the switch degree of an indirect network are analogous to the dimension and radix of a direct network. However, the bisection width of an indirect network is independent of switch degree. Because indirect networks do not exploit locality it is not possible to trade off diameter for bandwidth. When wire density is the limiting resource, a high-bandwidth direct network is preferable to an indirect network.

The low-dimensional k -ary n -cube provides a very general communication media for digital systems. These networks have been developed primarily for message-passing concurrent computers. They could also be used in place of a bus or indirect network in a shared-memory concurrent computer, in place of a bus to connect the components of a sequential computer, or to connect subsystems of a special purpose digital system. With VLSI communication chips, the cost of implementing a network node is comparable to the cost of interfacing to a shared bus, and the performance of the network is considerably greater than the performance of a bus.

The networks described here have been demonstrated in the laboratory and incorporated into commercial multiprocessors. The Torus Routing Chip (TRC) is a VLSI chip designed to implement low-dimensional k -ary n -cube interconnection networks [9]. The TRC performs wormhole routing in arbitrary k -ary n -cube interconnection networks. A single TRC provides 8-bit data channels in two dimensions and can be

cascaded to add more dimensions or wider data channels. A TRC network can deliver a 150-bit message in a 1024 node 32-ary 2-cube with an average latency of 7.5 μ s, an order of magnitude better performance than would be achieved by a binary n -cube with bit-serial channels. A new routing chip, the Network Design Frame (NDF), improves the latency to $\approx 1 \mu$ s [12]. The Ametek 2010 uses a 16-ary 2-cube (without end around connections) for its interconnection network [1].

Now that the latency of communication networks has been reduced to a few microseconds, the latency of the processing nodes T_{node} dominates the overall latency. To efficiently make use of a low-latency communication network we need a processing node that interprets messages with very little overhead. The design of such a *message-driven processor* is currently underway [7], [11].

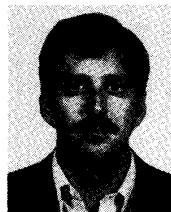
The real challenge in concurrent computing is software. The development of concurrent software is strongly influenced by available concurrent hardware. We hope that by providing machines with higher performance internode communication we will encourage concurrency to be exploited at a finer grain size in both system and application software.

ACKNOWLEDGMENT

I thank C. Seitz of Caltech for his many helpful suggestions during the early stages of this research.

REFERENCES

- [1] Ametek Corporation, Ametek 2010 product announcement, 1987.
- [2] K. E. Batcher, "Sorting networks and their applications," in *Proc. AFIPS FJCC*, vol. 32, 1968, pp. 307-314.
- [3] K. E. Batcher, "The Flip network in STARAN," in *Proc. 1976 Int. Conf. Parallel Processing*, pp. 65-71.
- [4] V. E. Benes, *Mathematical Theory of Connecting Networks and Telephone Traffic*. New York: Academic, 1965.
- [5] L. N. Bhuyan and D. P. Agrawal, "Generalized hypercube and hyperbus structures for a computer network," *IEEE Trans. Comput.*, vol. C-33, no. 4, pp. 323-333, Apr. 1984.
- [6] S. Browning, "The tree machine: A highly concurrent computing environment," Dep. Comput. Sci., California Instit. Technol., Rep. 3760, 1980.
- [7] W. J. Dally, *A VLSI Architecture for Concurrent Data Structures*. Hingham, MA: Kluwer, 1987.
- [8] W. J. Dally and C. L. Seitz, "Deadlock-free message routing in multiprocessor interconnection networks," *IEEE Trans. Comput.*, vol. C-36, no. 5, pp. 547-553, May 1987.
- [9] —, "The torus routing chip," *J. Distributed Syst.*, vol. 1, no. 3, pp. 187-196, 1986.
- [10] W. J. Dally, "Wire efficient VLSI multiprocessor communication networks," in *Proc. Stanford Conf. Advanced Res. VLSI*, Losleben, Ed. Cambridge, MA: MIT Press, Mar. 1987, pp. 391-415.
- [11] W. J. Dally *et al.*, "Architecture of a message-driven processor," in *Proc. 14th ACM/IEEE Symp. Comput. Architecture*, June 1987, pp. 189-196.
- [12] W. J. Dally and P. Song, "Design of a self-timed VLSI multicomputer communication controller," in *Proc. IEEE Int. Conf. Comput. Design*, 1987.
- [13] P. Kermani and L. Kleinrock, "Virtual cut-through: A new computer communication switching technique," *Comput. Networks*, vol. 3, pp. 267-286, 1979.
- [14] D. H. Lawrie, "Alignment and access of data in an array processor," *IEEE Trans. Comput.*, vol. C-24, no. 12, pp. 1145-1155, Dec. 1975.
- [15] C. L. Leiserson, "Fat trees: Universal networks for hardware-efficient supercomputing," *IEEE Trans. Comput.*, vol. C-34, no. 10, pp. 892-901, Oct. 1985.
- [16] C. A. Mead and L. A. Conway, *Introduction to VLSI Systems*. Reading, MA: Addison-Wesley, 1980.
- [17] M. C. Pease, III, "The indirect binary n -cube microprocessor array," *IEEE Trans. Comput.*, vol. C-26, no. 5, pp. 458-473, May 1977.
- [18] G. F. Pfister and V. A. Norton, "Hot spot contention and combining in multistage interconnection networks," *IEEE Trans. Comput.*, vol. C-34, no. 10, pp. 943-948, Oct. 1985.
- [19] C. L. Seitz, "Concurrent VLSI architectures," *IEEE Trans. Comput.*, vol. C-33, no. 12, pp. 1247-1265, Dec. 1984.
- [20] C. L. Seitz *et al.*, "The hypercube communications chip," Dep. Comput. Sci., California Inst. Technol., Display File 5182:DF:85, Mar. 1985.
- [21] C. H. Sequin, "Single chip computers, The new VLSI building block," in *Proc. Caltech Conf. VLSI*, C. L. Seitz, Ed., Jan. 1979, pp. 435-452.
- [22] H. J. Siegel, "Interconnection network for SIMD machines," *IEEE Comput. Mag.*, vol. 12, no. 6, pp. 57-65, June 1979.
- [23] H. S. Stone, "Parallel processing with the perfect shuffle," *IEEE Trans. Comput.*, vol. C-20, no. 2, pp. 153-161, Feb. 1971.
- [24] H. Sullivan and T. R. Bashkow, "A large scale homogeneous machine," in *Proc. 4th Ann. Symp. Comput. Architecture*, 1977, pp. 105-124.
- [25] A. S. Tanenbaum, *Computer Networks*. Englewood Cliffs, NJ: Prentice-Hall, 1981.
- [26] C. D. Thompson, "A complexity theory of VLSI," *Dep. Comput. Sci., Carnegie-Mellon Univ.*, Tech. Rep. CMU-CS-80-140, Aug. 1980.



William J. Dally (S'78-M'86) received the B.S. degree in electrical engineering from Virginia Polytechnic Institute, Blacksburg, in 1980, the M.S. degree in electrical engineering from Stanford University, Stanford, CA, in 1981, and the Ph.D. degree in computer science from Caltech in 1986.

From 1980 to 1982, he worked at Bell Telephone Laboratories where he contributed to the design of the BELLMAC-32 microprocessor. From 1982 to 1983 he worked as a consultant in the area of digital systems design. From 1983 to 1986 he was a Research Assistant and then a Research Fellow at Caltech. He is currently an Associate Professor of Computer Science at the Massachusetts Institute of Technology. His research interests include concurrent computing, computer architecture, computer-aided design, and VLSI design.