

S.Q. Zheng  
Department of Computer Science  
Erik Jonsson School of Engineering and Computer Science  
University of Texas at Dallas  
Box 830688, MS EC 31  
Richardson, TX 75083-0688, USA  
Email: [sizheeng@utdallas.edu](mailto:sizheeng@utdallas.edu)

# Scalable and Practical Nonblocking Switching Networks

Part of the presented work was jointly carried out with Dr. Ashwin Guimare and Dr. Hong Shen.

---

## Acknowledgment

- Switching fabric in network switches and routers.
- Communication subsystem (interconnection network) in a parallel computing system.

---

## Applications of Switching Networks

- Multiple-multicasting Switching Networks.
- Permutation Switching Networks.

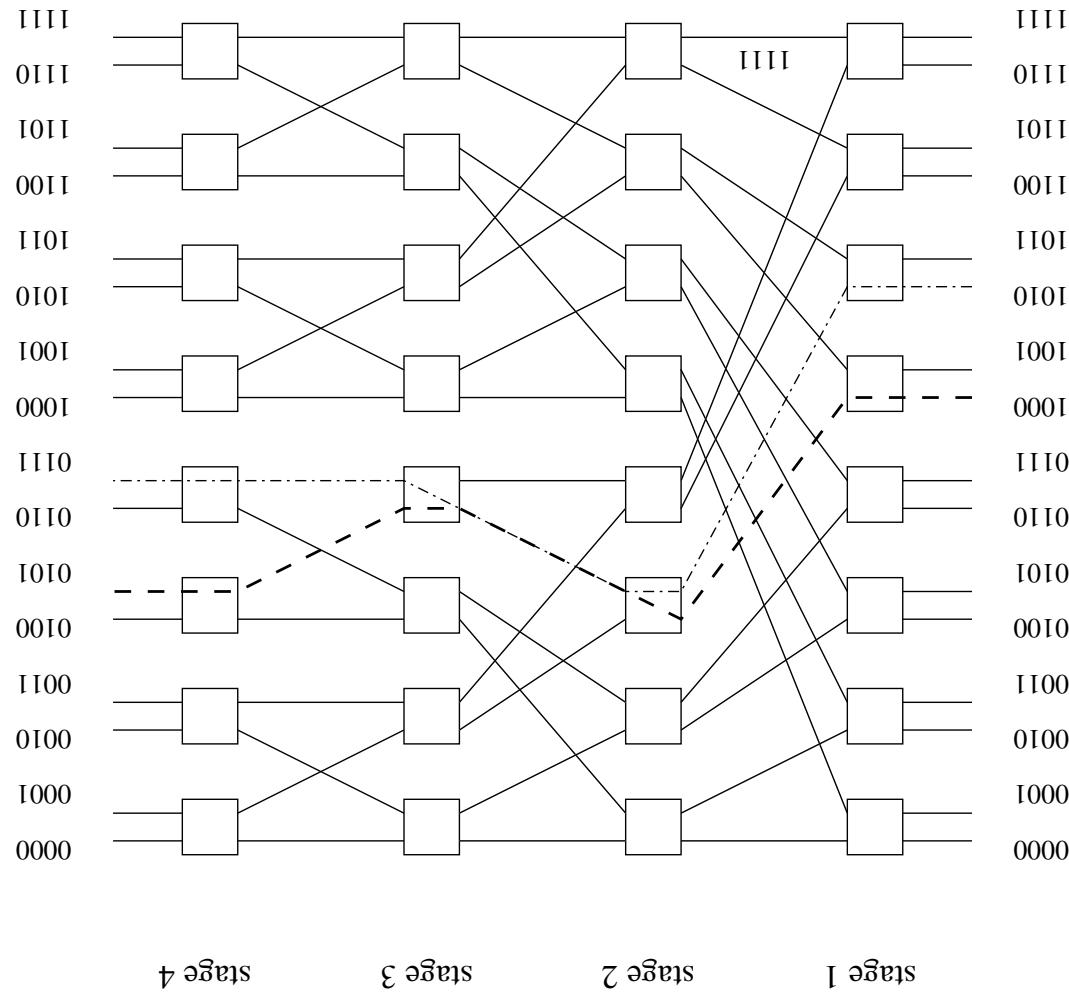
---

## Switching Networks

An  $N \times N$  switching network has  $N$  inputs and  $N$  outputs. It usually comprises a number of switching elements (SEs) interconnected by a set of links. Without loss of generality, we assume that an SE is of size  $2 \times 2$ ; i.e., the two inputs of an SE intend to be connected with the same output causes link conflict. If an I/O connection path does not have any link conflict with other connection paths, it is called a conflict-free path.

Nonblocking switching networks have been favored in switching systems because they can be used to set up conflict-free I/O connection paths for any (partial) I/O permutation.

## A baseline network.



**Example**

- **Three Basic Types of Nonblocking Switching Networks**
- A connection from any idle input to any idle output can always be established regardless of how existing connections were established.
  - **Wide-Sense Nonblocking (WSNB) networks**
  - Useful for both circuit switching and packet switching.
  - A connection from any idle input to any idle output can always be established by following a rule (algorithm).
  - Useful for both circuit switching and packet switching.
- A connection from any idle input to any idle output can always be established.
- **Strictly Nonblocking (SNB) networks**.

- Rearrangeable Nonblocking (RNB) networks.
  - A connection from any idle input to any idle output can be established if rearrangement of existing connections is allowed.
  - Useful for packet (cell) switching but not suitable for circuit switching.

$O(N \log N)$  cost practically achievable. Example: Benes network.

---

## Optimal RNB Network

Suboptimal: Clos network of  $O(N^{1.5})$  cost, and Cantor network of  $O(N \log_2 N)$  cost. Routing I/O connections is time-consuming.

$O(N \log N)$ -cost theoretically possible, but not practical.

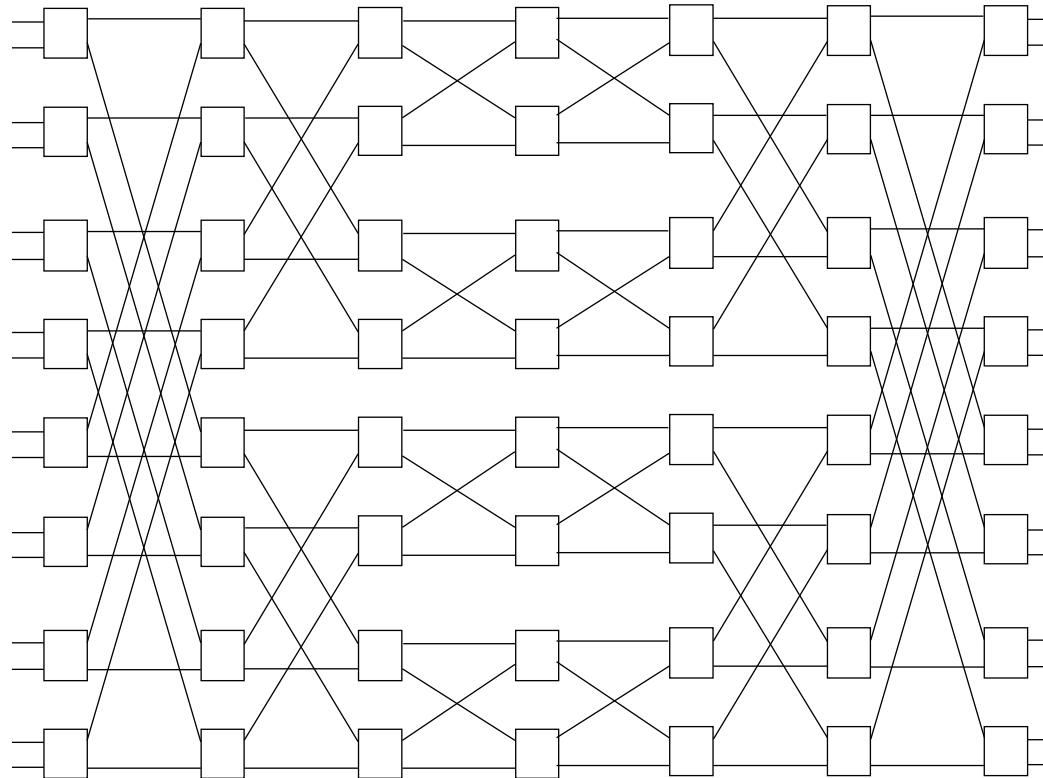
---

## Optimal SNB and WSNB Networks

Cost:  $\mathcal{O}(N \log N)$  SESs.

---

Lower Bound for the Cost of Nonblocking Switching Networks works



---

**Benes Network**

Practically, maybe not.

---

Question: Do we really need SNB and WSNB networks?

Find practical  $O(N \log N)$ -cost SNB and WSNB networks for circuit-switching applications.

---

An Outstanding Open Problem

This definition is broad enough to cover many embodiments of VNB network architectures.

(iii) when properly controlled, the network can perfectly emulate an WSNB network.

(i) it provides more than one path from any input (i.e. source) to any output (i.e. destination) such that different paths that connect the same source/destination pair can be momentarily set up simultaneously but any two paths connecting different input/output pairs are link-disjoint at any time; and

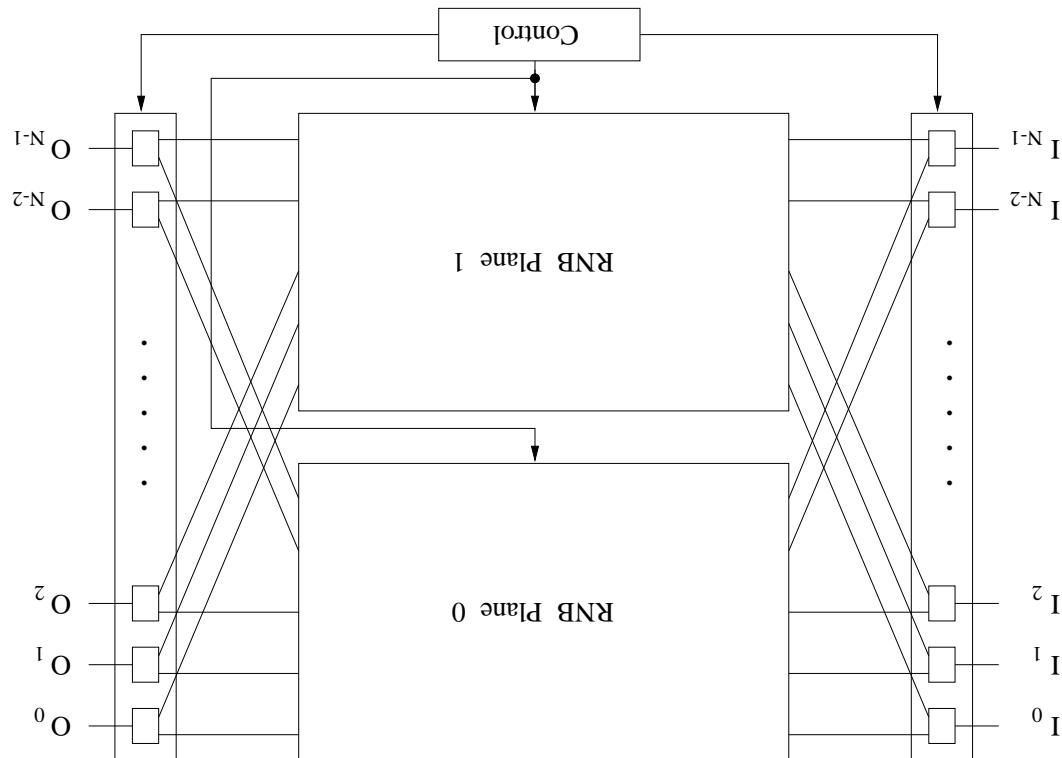
flies the following conditions:

A switching network is a virtual nonblocking (VNB) network if it satisfies

---

## Virtual Nonblocking Networks

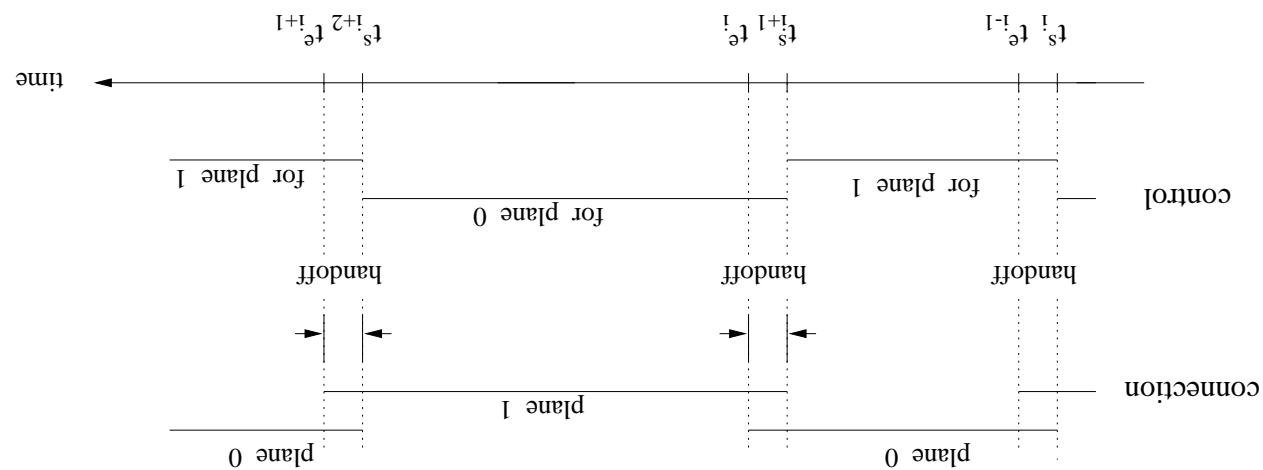
## Ping-Pong VNB network architecture.



---

## Ping-Pong VNB Network

## Timing diagram of the operations of a ping-pong VNB network.



## Handoff Process of Ping-Pong VNB Network

## Handoff Schemes

Lazy handoff scheme: handoff is only needed when finding conflict-free paths fails.

Eager handoff scheme: handoff is invoked when new connection requests arrive, regardless of whether conflict-free paths for the new connection requests exist or not.

VNB networks are useful in circuit switched networks. They are used as switching networks within network routers and switches to establish and maintain nondisruptive circuits.

## Self-routing VNB networks are desirable.

Generalization: This property can be generalized to include the case that the setup of each SE only depends on the  $O(\log N)$  information, which may not be the addresses of source and destination, available at the SE during routing.

Self-routing: the ability that any connection in a network can be established by inspecting addresses of its source and destination regardless of other connections.

Parallel Benes routing algorithm using an  $N$ -processor hypercube takes  $O(\log_4 N)$  time, which is too slow.

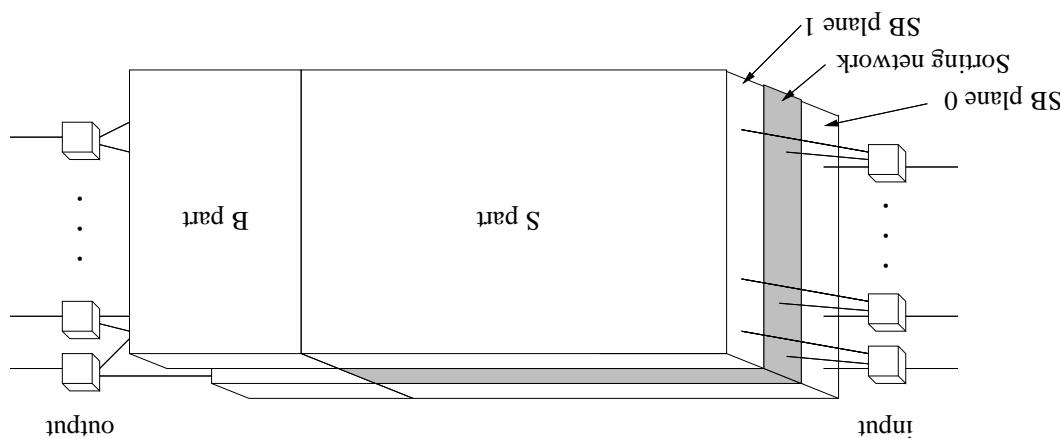
## Routing Performance Consideration

The ping-pong VNB network constructed using two copies of Benes network has cost  $O(N \log N)$ , which is optimal.

## Optimal VNB Networks

Using AKS sorting network, this VNB network has cost  $O(N \log N)$ . But this is not practical because the coefficient hidden by big-O is too large.

An SB ping-pong VNB network.



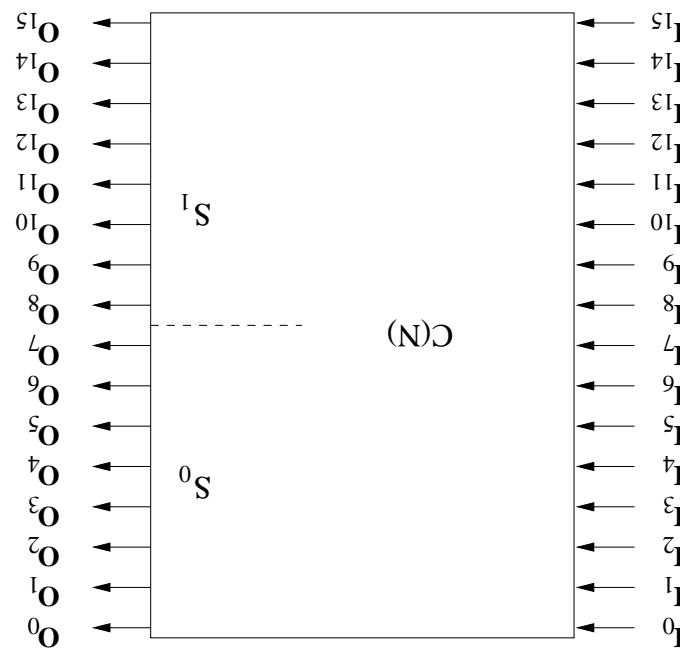
Each plane is a concatenation of a Sorting network and a Banyan (but terribly, Omega, baseline) network.

---

**SB VNB Networks**

$C(N)$  is an  $N \times N$  network that separates a set  $S$  of  $N$  input numbers into two subsets  $S_0$  and  $S_1$  such that  $|S_0| = |S_1| = \frac{N}{2}$ , no number in  $S_0$  is larger than any number in  $S_1$ , and numbers in  $S_0$  and  $S_1$  appear at the first and second half of the outputs, respectively.

Classifier  $C(16)$ .

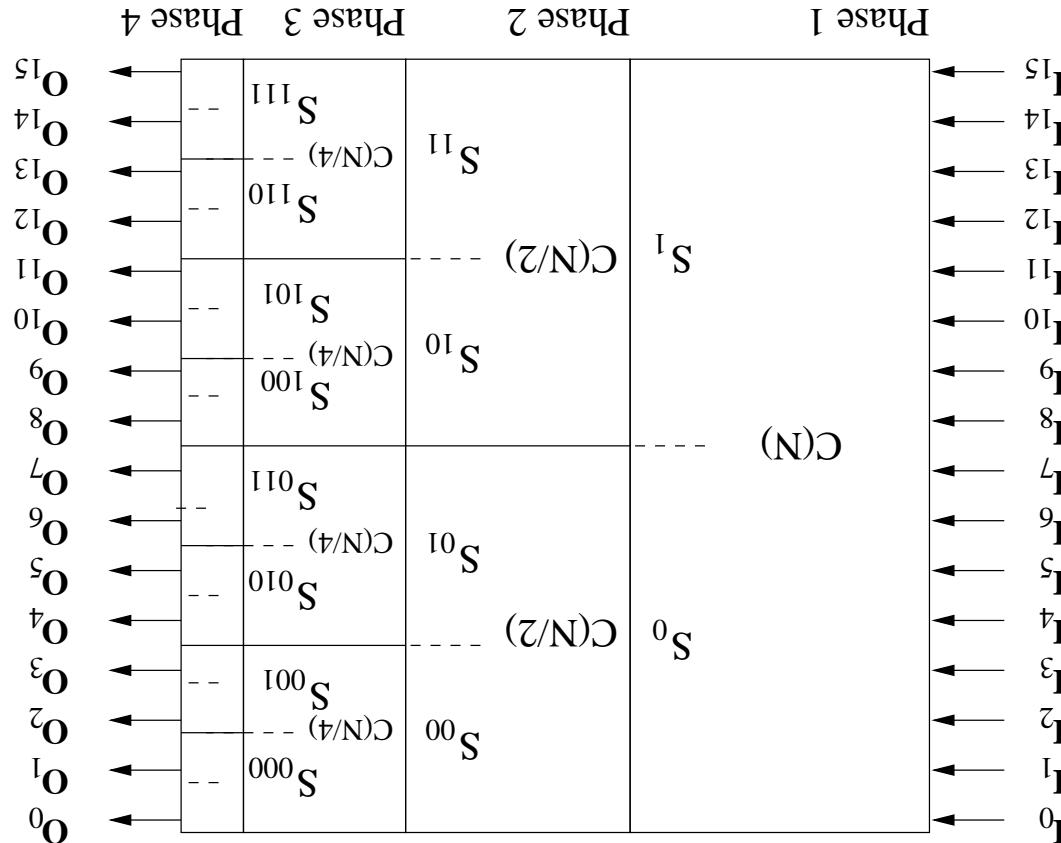


Building block:  $N$ -classifier  $C(N)$ .

---

Q-Sort: A New Sorting Network

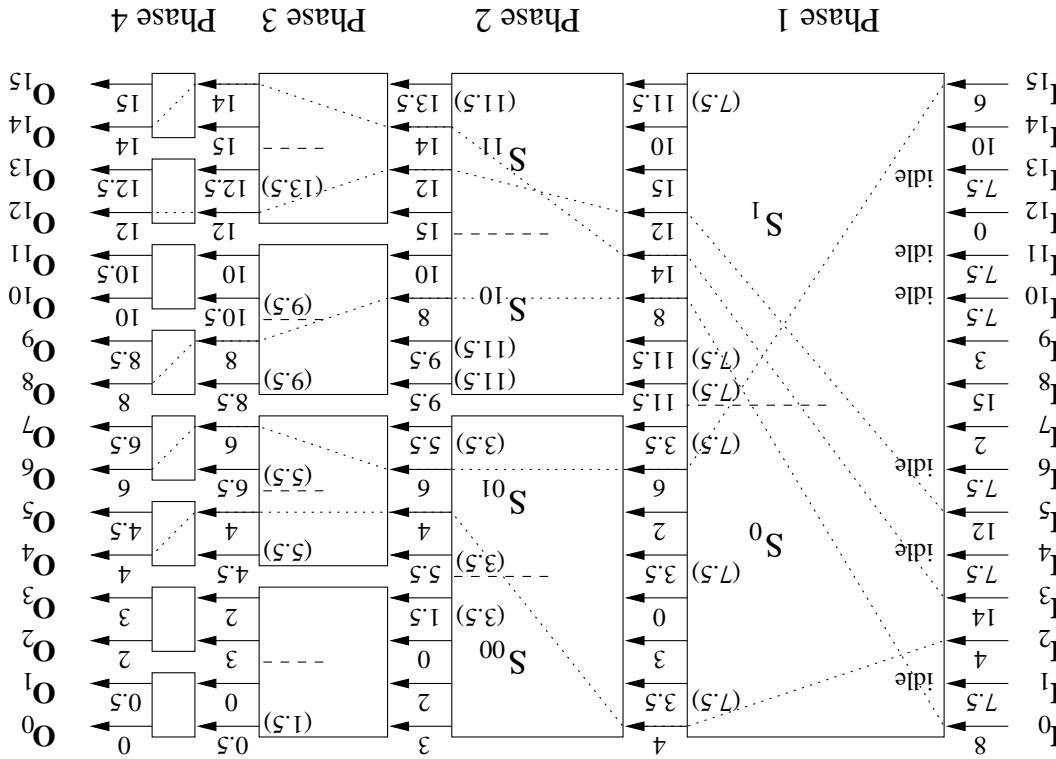
A  $\mathcal{Q}$ -sort network  $\mathcal{Q}(N)$ ,  $N = 16$ . Cost:  $O(N \log_2 N)$ .



**$\mathcal{Q}$ -Sort Network  $\mathcal{Q}(N)$ : Constructed Using Classifiers Re-**

curatively

**Self-routing process in  $Q(N)$  for partial permutation  $\pi$ , where  $h_i = \forall$  indicates that input  $i$  is idle (not to be connected to output).**



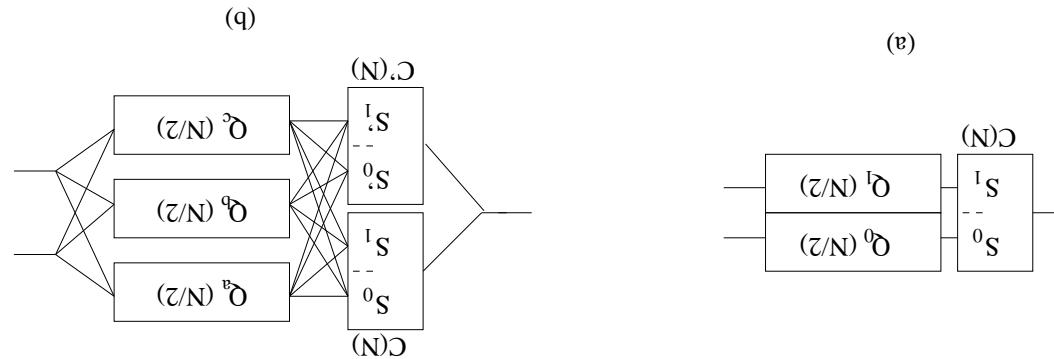
$$\pi = (h_0, h_1, \dots, h_{N-1}) = (8, A, 4, 14, A, 12, A, 2, 15, 3, A, A, 0, A, 10, 6)$$

**Routing a partial permutation**

---

**$Q$ -Sort Network  $Q(N)$  as an RNB Network**

(a) The structures of a  $Q$ -sort network  $\tilde{Q}(N)$ . (b) The cost-saving VNB network  $\tilde{Q}_*(N)$  based on  $\tilde{Q}(N)$ .



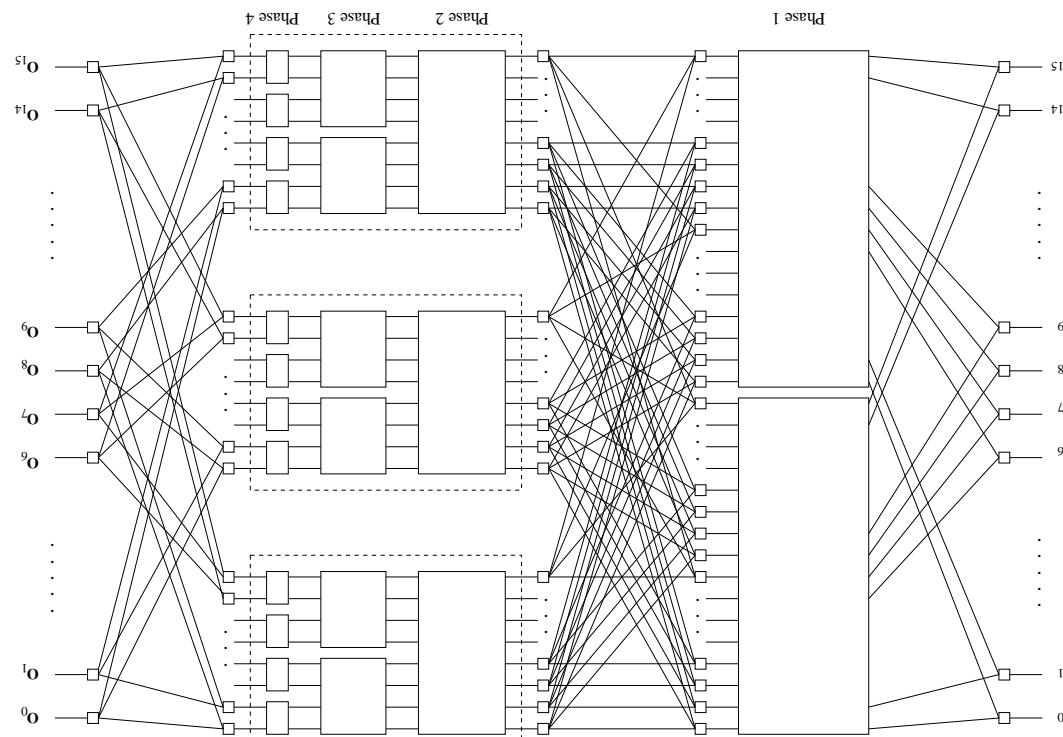
## $Q_*(N)$ Network: Cost Saving VNB Network

Using two copies of  $Q$ -sort network, we can construct a self-routing ping-pong VNB network.

## Self-Routing Ping-Pong VNB Network Based in $Q(N)$

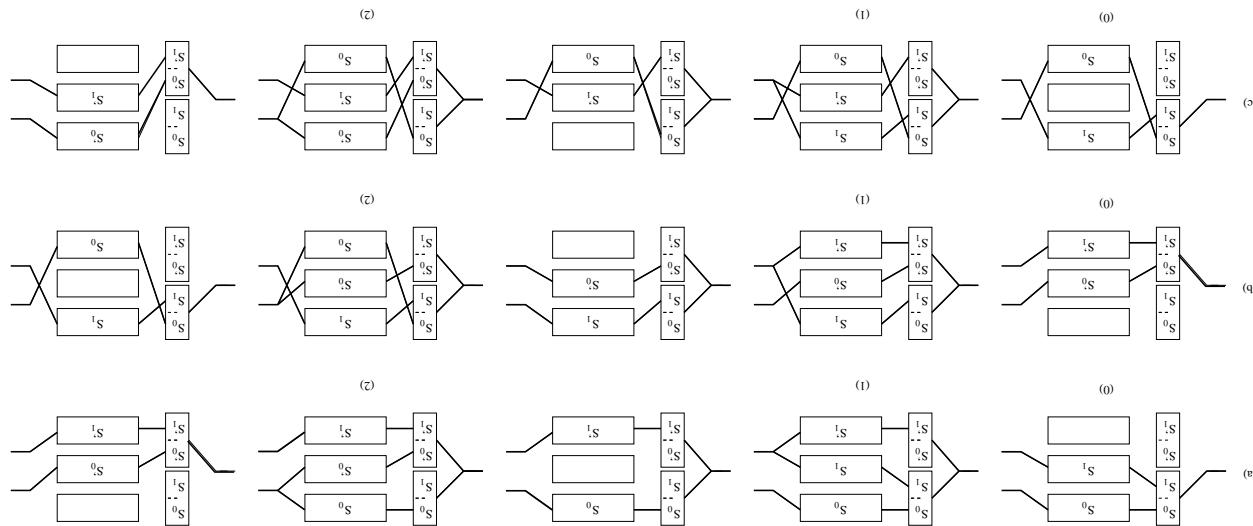
$Q(8)$ .

The interconnection details of  $Q_*(16)$ . The part in a dashed block is a



More Details of  $Q_*(N)$  Network

## Handoff operations of $\mathcal{Q}_*(N)$ .



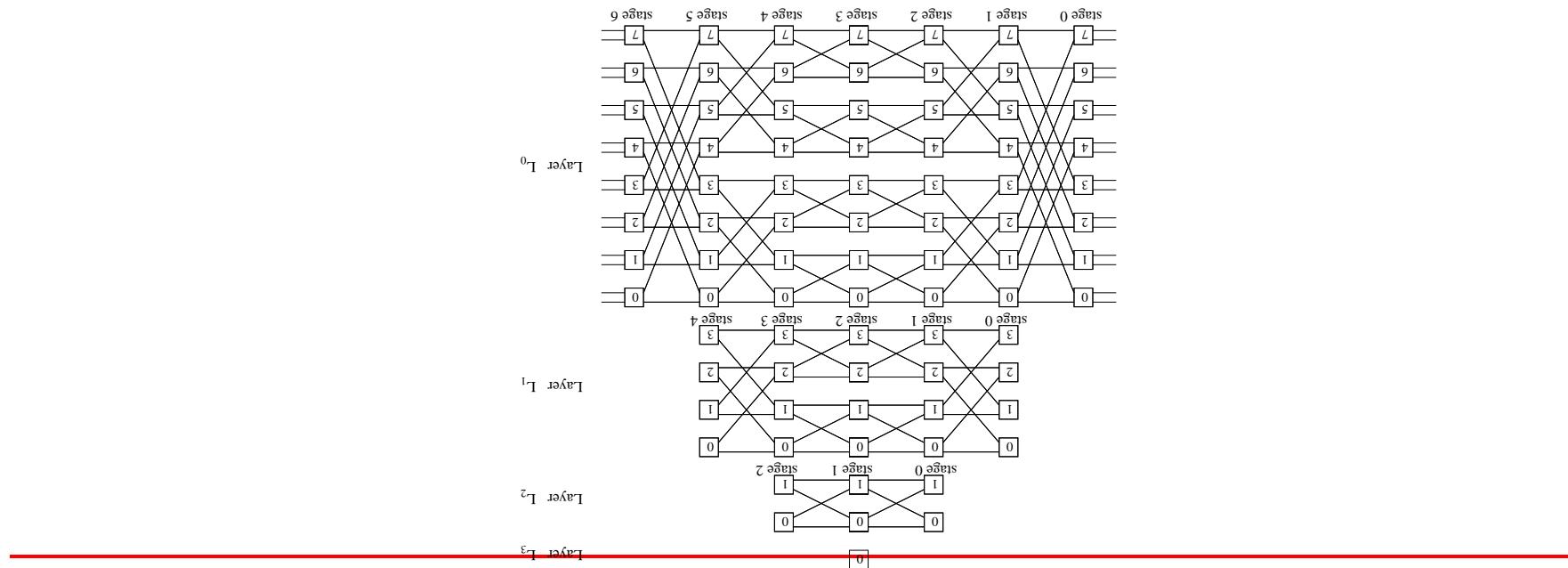

---

## Cyclic Handoff Scheme for $\mathcal{Q}_*(N)$ Network

- The concept of virtual nonblocking (VNB) networks, which are operationally equivalent to SNB or WSNB networks but as simple as RNB networks.
- The ping-pong structure, for designing a class of VNB networks.
- It is simple to design an optimal VNB network that contains  $O(N \log N)$  SESes using two Benes networks. In contrast, no explicitly constructed SNB network of cost in the order of  $O(N \log N)$  cost with a small coefficient has been reported in the literature.
- Theoretically it is possible to construct self-routing VNB networks with  $O(N \log N)$  SESs. Practically it is easy to construct self-routing VNB networks with  $O(N \log^2 N)$  SESs by adopting some routing strategy. To the best of our knowledge, no self-routing SNB network of  $O(N \log^2 N)$  SESs has been reported in the literature.

- A new sorting network, called the  $\mathcal{Q}$ -sort network, based on classifiers, and show how to construct self-routing RNB networks with  $O(N \log^2 N)$  cost.
- A cyclic handoff scheme for constructing VNB network  $\mathcal{Q}_*(N)$  with reduced cost.
- VNB networks can achieve lower cost and faster routing, which are not easy to achieve simultaneously for SNB networks. Thus, for all circuit switching applications, all we need are VNB networks.

## Topology of $PB(16)$ : intra-layer connections.

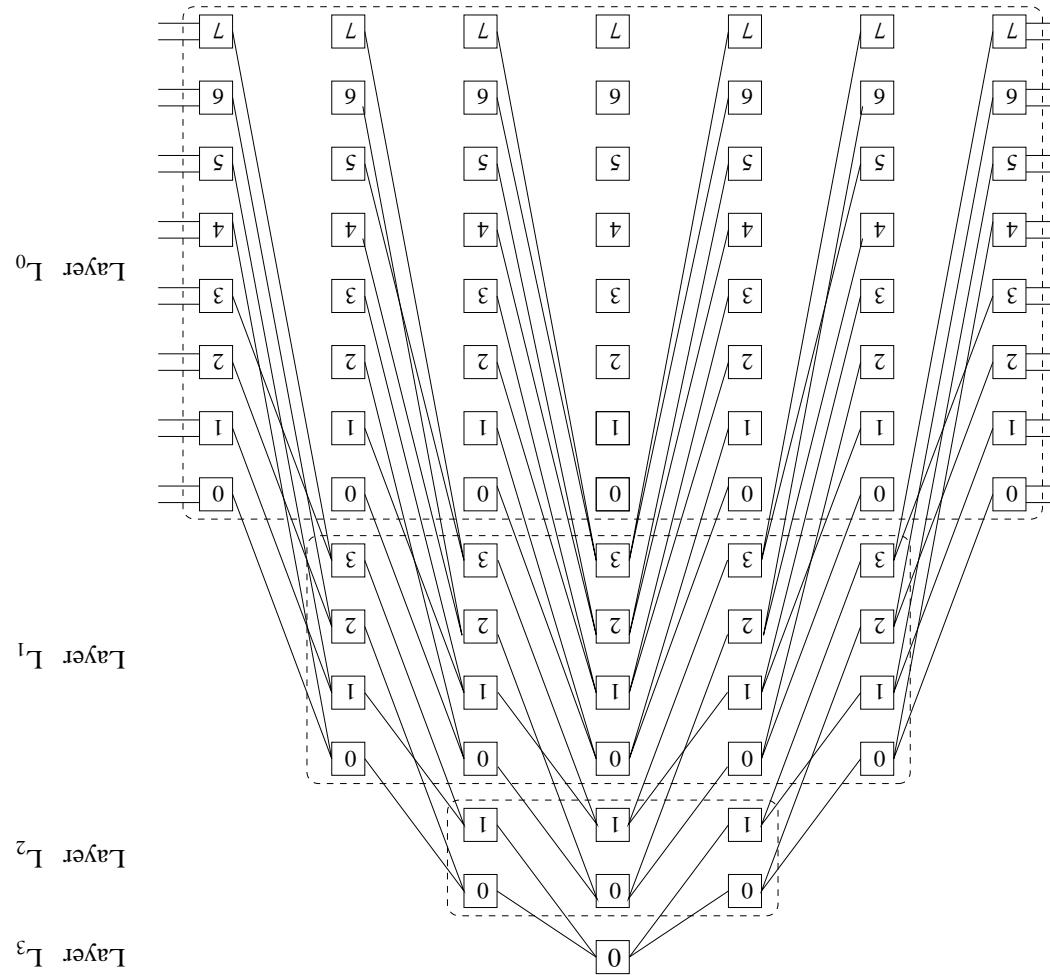


## Intra-layer Connections

Pyramid Benes network  $PB(2^n)$  has  $n - 1$  layers  $L^i$ ,  $0 \leq i \leq n - 2$ , such that  $L^i$  is a  $2^{n-i} \times 2^{n-i}$  Benes network.

## Almost-Nonblocking Permutation Switching Networks

## Topology of PB(16): Inter-layer connections.



Inter-layer Connections

## Nonblockingness of $PB(N)$ Networks

$PB(N)$  has  $O(N \log N)$  cost.

By extensive simulations, we have not found any blocking case for  $PB(N)$ .

We cannot claim that  $PB(N)$  is strictly nonblocking.

However, we can claim that  $PB(N)$  is almost strictly nonblocking.

Broadcasting is a communication pattern which sends a message from a given input port  $i$  to all output ports.

Multicast is a communication pattern that sends a message to a given input port  $i$  to a subset  $O^i$  of output ports.

A multicast I/O connection pattern is represented by a pair  $(i, O^i)$ .

A legal multiple-multicast is a communication pattern that is represented by a set of pairs  $\{(i_1, O^{i_1}), (i_1, O^{i_1}), \dots, (i_m, O^{i_m})\}$  such that  $i_j \neq i_k$  and  $O^{i_j} \cup O^{i_k} = \emptyset$  if  $j \neq k$ .

Permutation (i.e one-to-one) and broadcast (i.e. one-to-all) are special cases of legal multiple-multicast.

In order to obtain a VNB network with multiple-multicast capability,  
we have to design a RNB network with multiple-multicast capability.

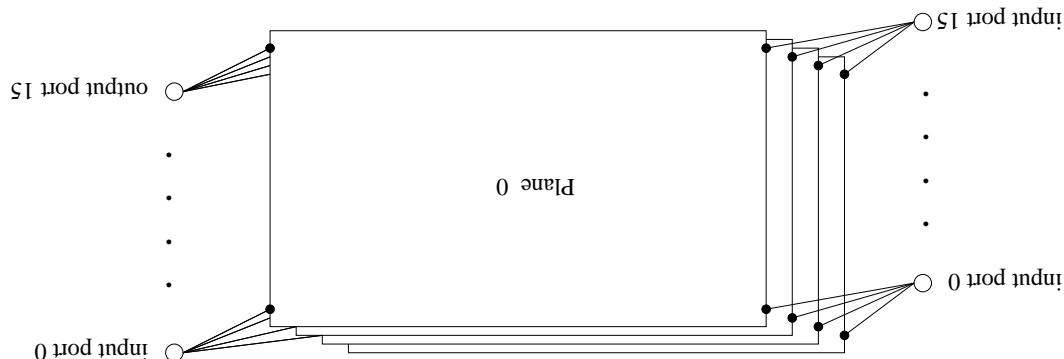
---

## VNB Networks for Multiple-Multicast

No known multi- $\log_2 N$  of cost  $O(N^{1.5} \log N)$  can realize conflict-free multiple-multicast.

If each plane is a Banyan network, then  $p \geq 2 \lfloor \frac{n}{2} \rfloor$  planes are required for the multi- $\log_2 N$  network to be an RNB permutation network. Cost:  $O(N^{1.5} \log N)$ .

**Multi- $\log_2 N$  networks.**

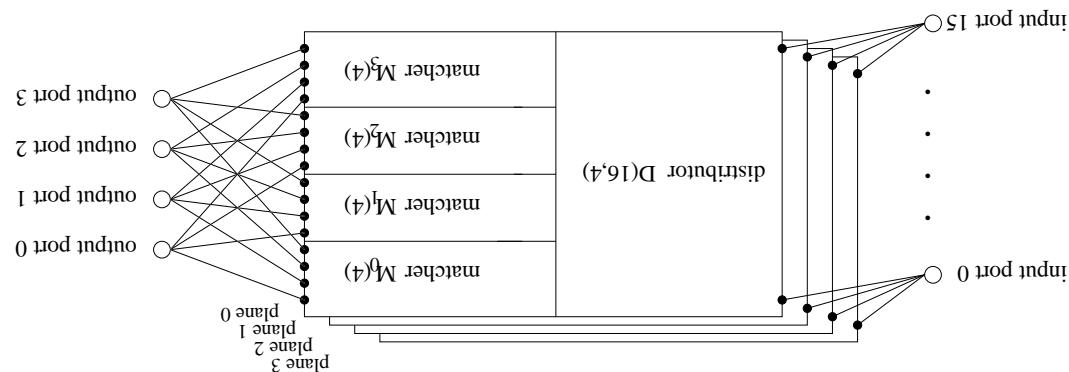


An  $N \times N$  network of  $p$  identical planes, each being a multistage net-work of  $O(\log_2 N)$  stages.

**Multi- $\log_2 N$  Networks**

## $DM(N)$ : A Self-Routing Multi- $\log_2 N$ Network for Multiple-Multicast

### Structure of $DM(N)$ networks.



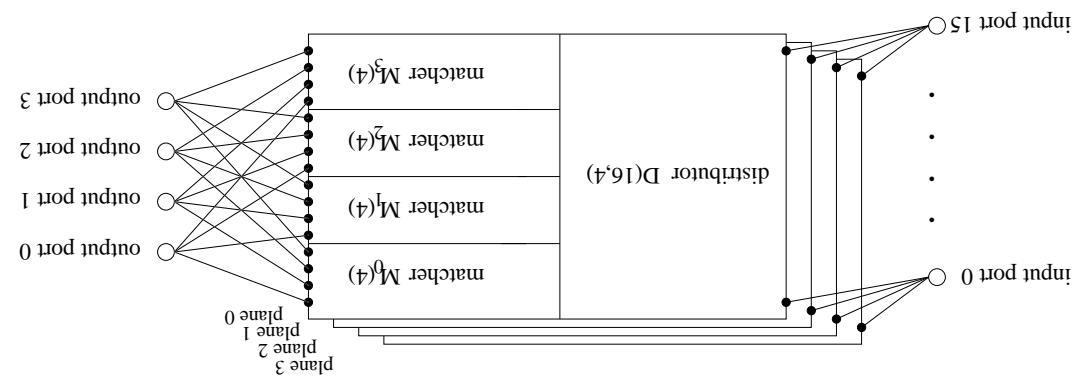
Assume that  $N = 2^n$ , where  $n$  is even.

---

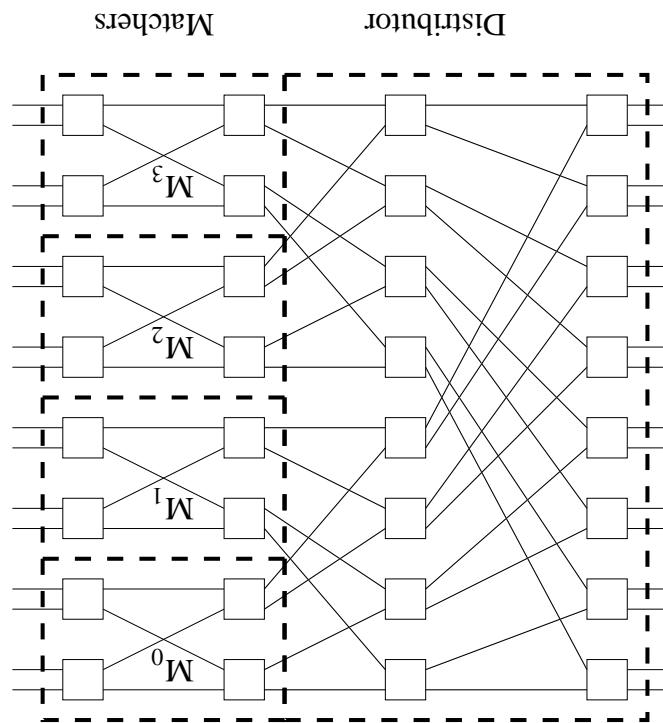
## Structure of $DM(N)$ Network

- $\sqrt{N}$  planes,  $P_0, P_1, \dots, P^{\sqrt{N}-1}$ , each being a baseline network.
- Input port  $i$  is connected to the  $i$ -th inputs of all planes through a tree  $T_i$ .
- Plane  $P_j$  is dedicated to connecting any input port to any output ports  $j \cdot \sqrt{N} + k$ ,  $0 \leq k < \sqrt{N}$ .
- Plane  $P_j$  has a distributor, and  $\sqrt{N}$  matchers, each having  $\sqrt{N}$  outputs, with the  $k$ th outputs of the matchers connected to output port  $j \cdot \sqrt{N} + k$ .

## Structure of $DM(N)$ networks.



A plane of  $DM(16)$  is a Baseline network.



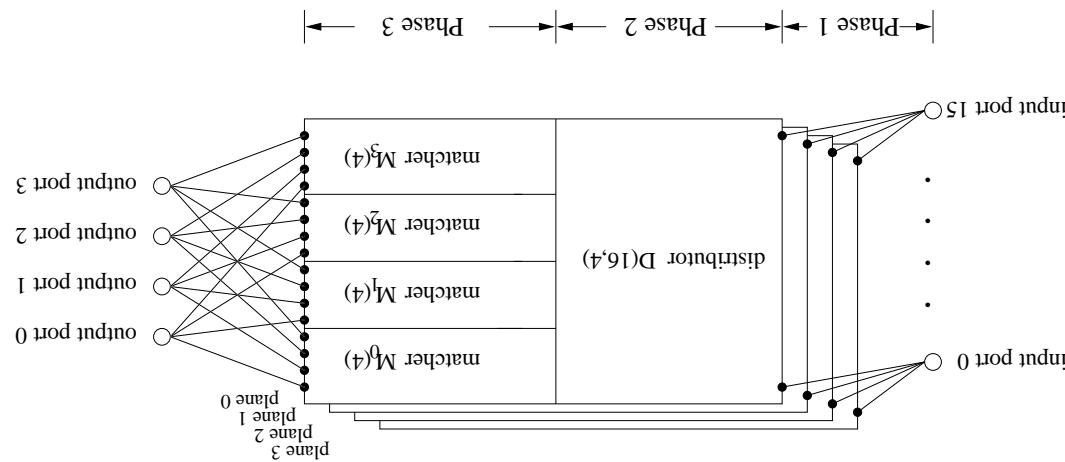
Structure of a Plane

An input is active if it is to be connected to an output.

---

### 3-Phase Routing Algorithm

3-phase routing in  $DM(N)$  networks.



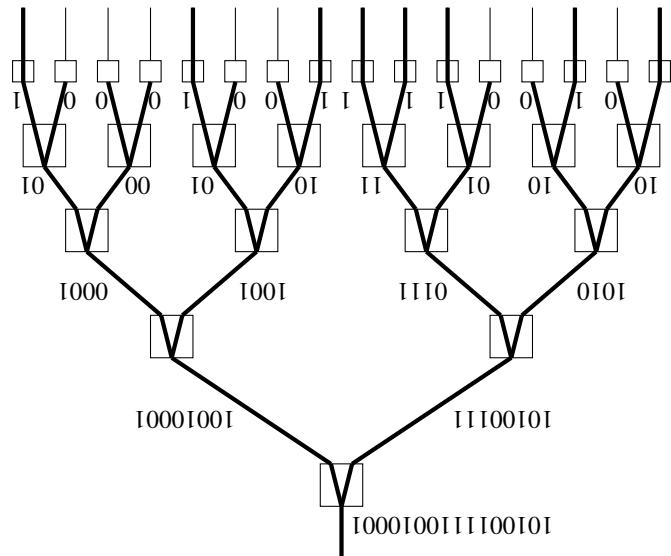
tree

Connect input port  $i_j$  to the  $k$ -th plane if  $O(i_j, k) \neq \emptyset$  using input

For all  $(i_j, O_{i_j})$  do the following in parallel:

**Phase 1: Route inputs to the inputs of planes.**

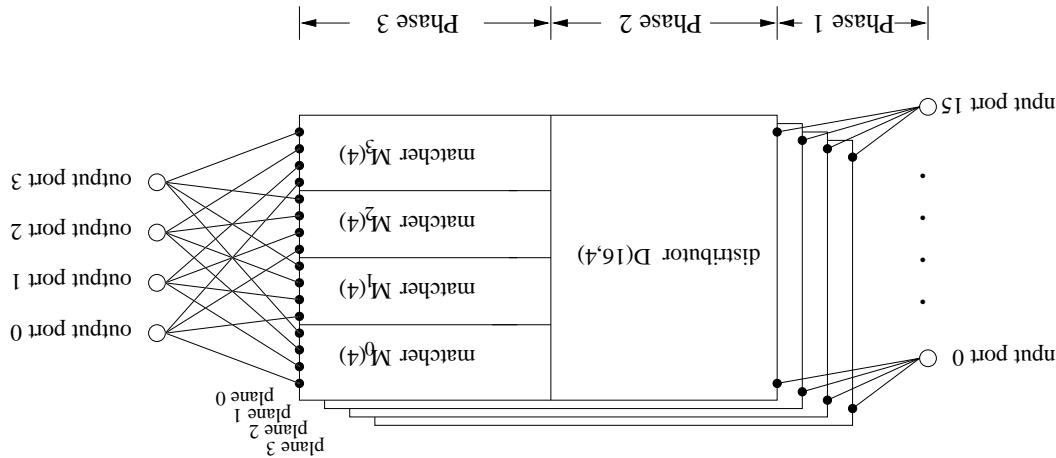
Broadcast-and-select: with a bit-string  $10100111001000$ , the input of the tree is connected to the outputs specified by the bit-string.



Each input port  $i_j$  computes a bit string for planes it has to be connected. Then, this bit string is used to route through the tree  $T_{i,j}$  using broadcast-and-select technique.

## Implementation of Phase 1

## 3-phase routing in $DM(N)$ networks.



Connect active inputs of each plane to the inputs of the plane's matchers such that each matcher has at most 2 active inputs.

For all planes do the following in parallel:

Phase 2: Routing through distributors.

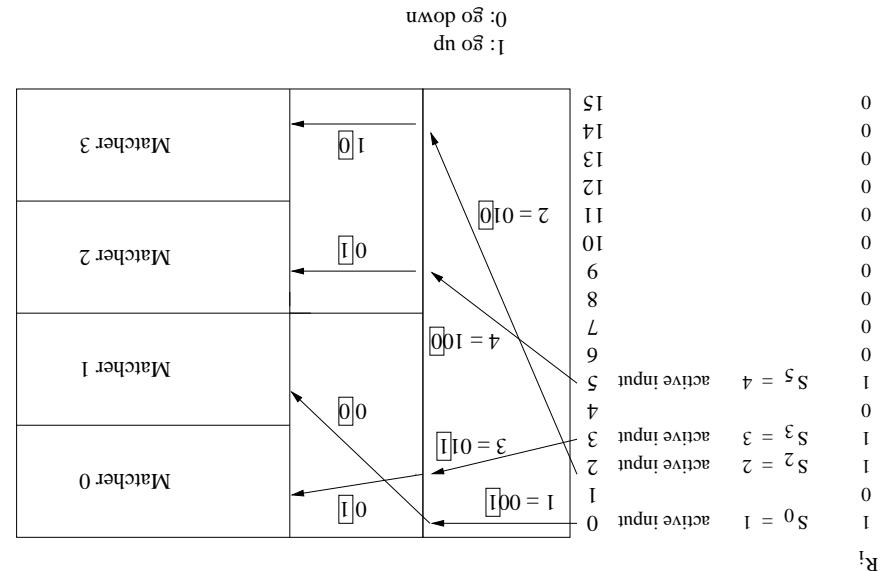
Each plane is equipped with a circuit for generating prefix sums of  $N$  Boolean values, which are used for self-routing in the plane. We call this circuit *intraplane routing preprocessing (IRP)* circuit. It is sufficient to consider routing operations in a single plane. If input  $i$  is active, set  $R_i = 1$ ; else set  $R_i = 0$ . The IRP circuit is used to compute binary prefix sums of  $R_i$ 's, where the binary prefix sum  $S_i$  is defined as

$$(1) \quad S_i = R_0 + R_1 + \cdots + R_i$$

$S_i$ 's can be computed in  $O(\log N)$  time by an IRP circuit of tree structure.

## Implementation of Phase 2

## Self-routing in Phase 2.

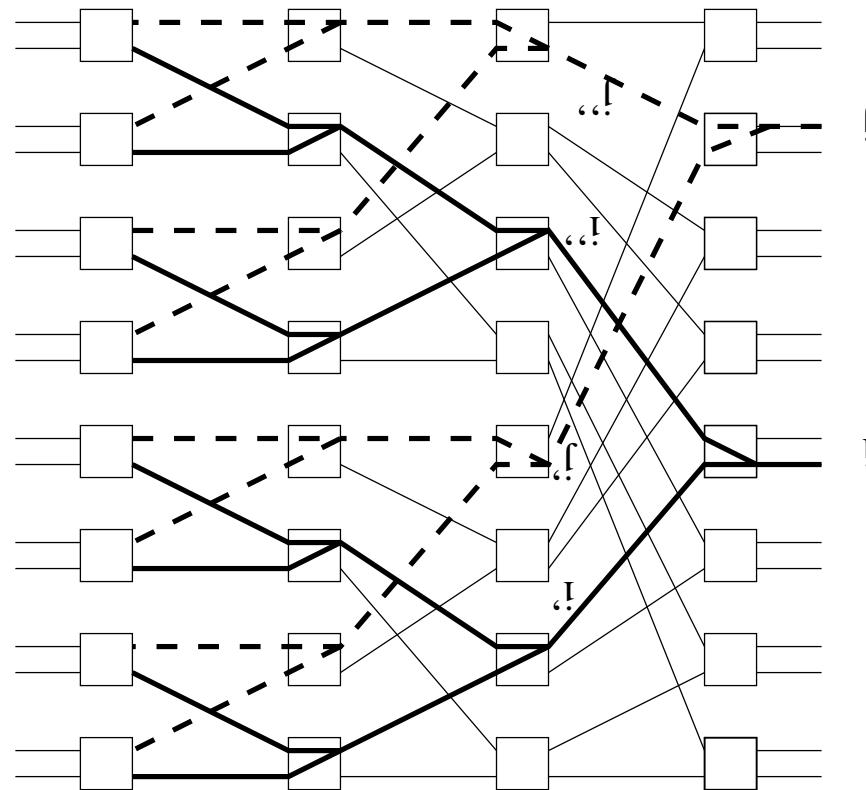


For  $R_0 = 1, R_1 = 0, R_2 = 1, R_3 = 0, R_4 = 1, R_5 = 1, R_6 = R_7 = R_8 = R_9 = R_{10} = R_{11} = R_{12} = R_{13} = R_{14} = R_{15} = 0$ , we have  $S_0 = 1, S_1 = 1, S_2 = 2, S_3 = 3, S_4 = 3, S_5 = 4, S_6 = 4, S_7 = 4, S_8 = 4, S_9 = 4, S_{10} = 4, S_{11} = 4, S_{12} = 4, S_{13} = 4, S_{14} = 4, S_{15} = 4$ .

---

**Example:**

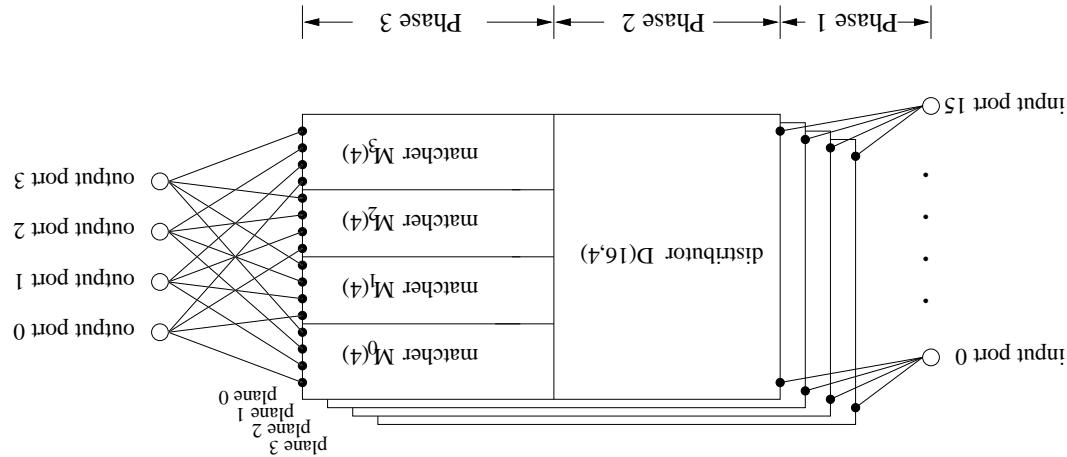
Two broadcasting trees in  $M(16)$ .



Phase 2 guarantees that each matcher has at most 2 active inputs.  
 Phase 2 guarantees the two active inputs have link-disjoint broadcast trees in the matcher, which is a Baseline subnetwork.

## Active Inputs of Matchers

### 3-phase routing in $DM(N)$ networks.



If an active input is originated from input port  $i_j$ , then connect the active input to all outputs of  $O_{i_j,k}$  within the matcher.

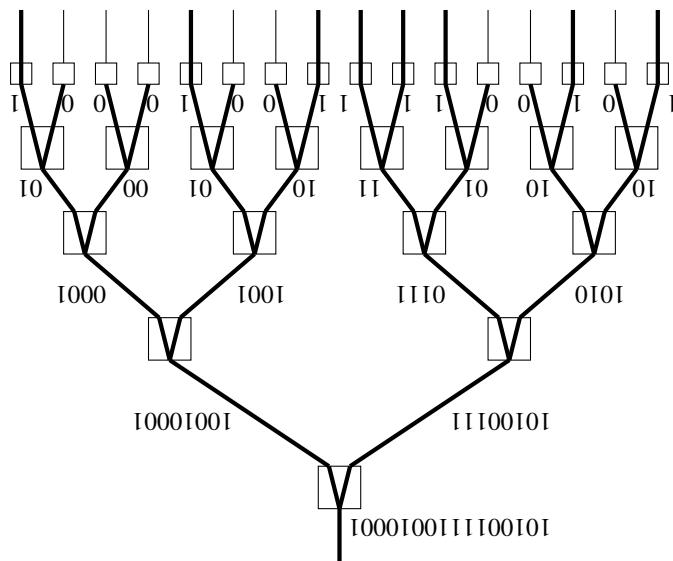
parallel:

For all matchers in all planes  $P_k$ ,  $0 \leq k < \sqrt{N}$ , do the following in

---

Phase 3: Routing through matchers.

## Broadcast-and-select.



Using broadcast-and-select technique based on bit strings..

---

Implementation of Phase 3

Then, self-routing is conducted in all phases.

of each plane.

Binary prefix sums for Phase 2 are performed by a circuit at the input

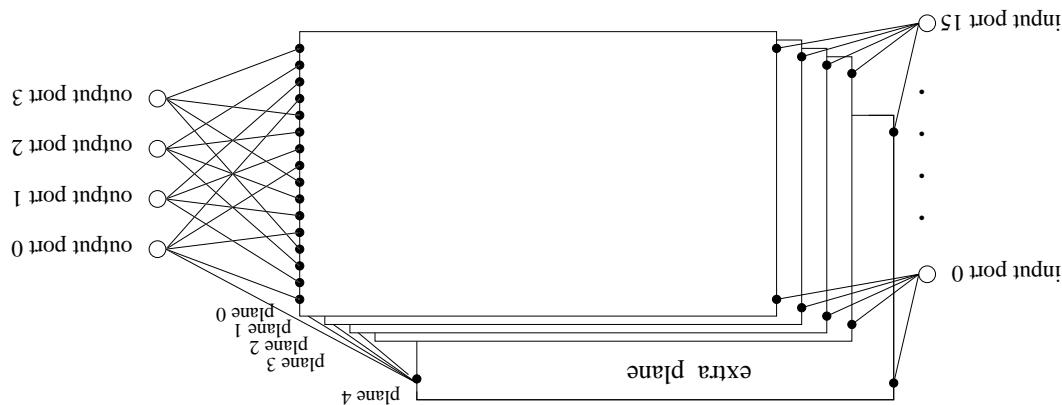
Bit strings for Phase 1 and Phase 3 are computed by input ports.

Routing in  $DM(N)$  is performed in 3 phases.

---

**Summary on Routing  $DM(N)$  Network**

## Extended $DM(16)$ network.



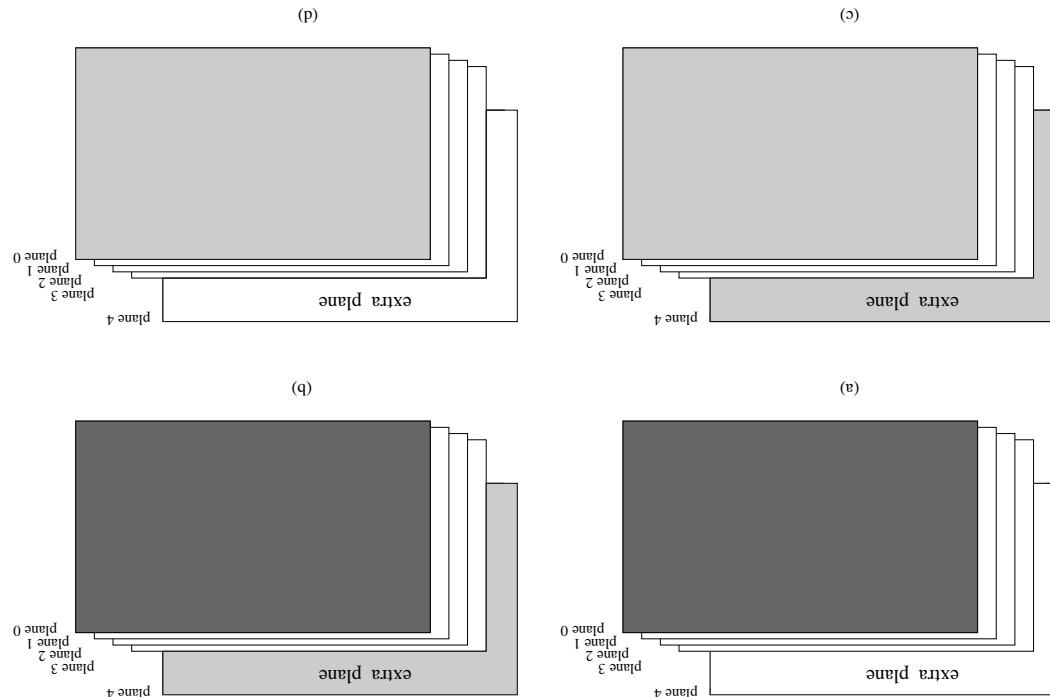
But we can reduce the cost by only introduce 1 extra plane.

Using ping-pong scheme, we need two copies of  $DM(N)$ .

---

Making  $DM(N)$  Virtually Nonblocking

## Using 1 extra plane to do hand-off.




---

**Hand-off in Virtually Nonblocking Multiple-Multicast  $DM(N)$**

**Network**

The high degree of connection capability in SNB and WSNB networks is achieved at much higher cost, which implies their lower scalability. We introduce the concept of virtual nonblockingness. We show that a virtual nonblocking network functions like a strictly or wide-sense nonblocking network, but it is constructed with the cost of a rearrangeable nonblocking network.

While finding low cost SNB networks remains to be a challenging theoretical and scientific pursuit, VNB networks can satisfy all the requirements of an SNB or WSNB in practice.

Our results indicate that for large-scale circuit switching applications, we only need to build virtual nonblocking networks.