

# **BANYAN-BASED SWITCHES**

**HIGH PERFORMANCE  
SWITCHES AND ROUTERS**

**Wiley**

**H. JONATHAN CHAO and BIN LIU**

**Instructor: Mansour Roustazadeh**

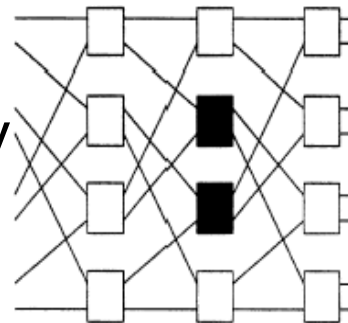
# INTRODUCTION

## Multistage Network

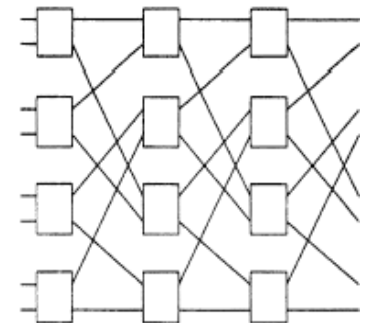
- First in circuit switched telephone networks
- To aim nonblocking with less crosspoints than a crossbar

## Banyan network

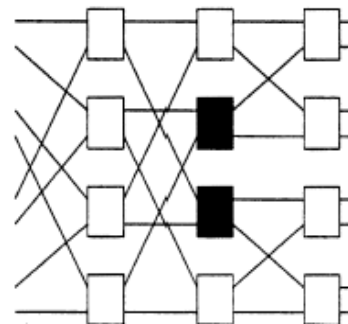
- Exactly one path from any input to any output
- 4 classes
  - a) Shuffle-exchange (Omega)
  - b) Reverse shuffle-exchange
  - c) Narrow-sense banyan
  - d) Baseline



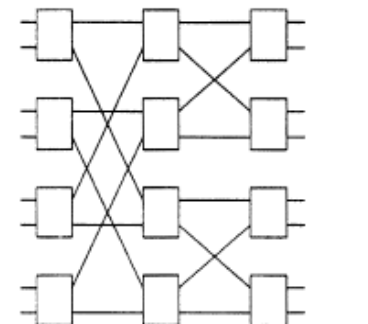
(a)



(b)



(c)

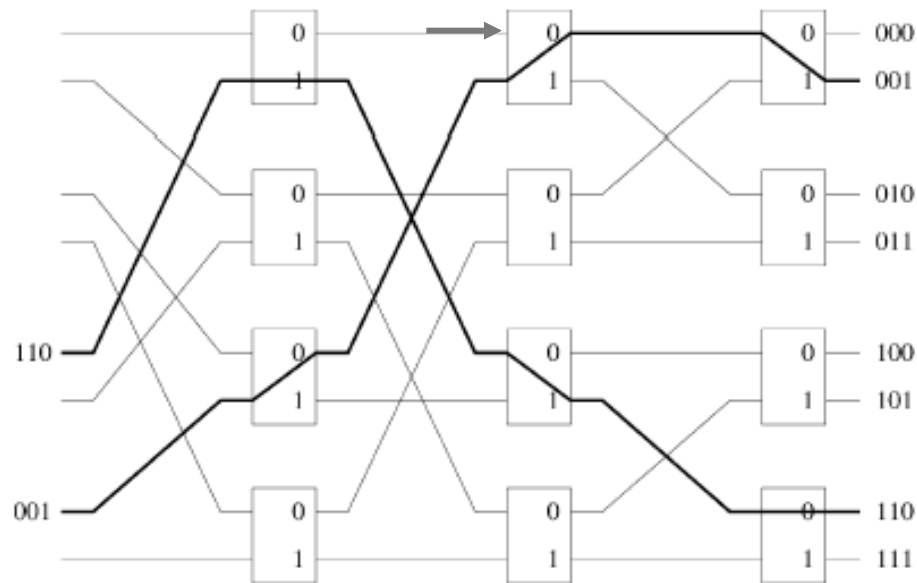


(d)

# INTRODUCTION

## Banyan network principal properties

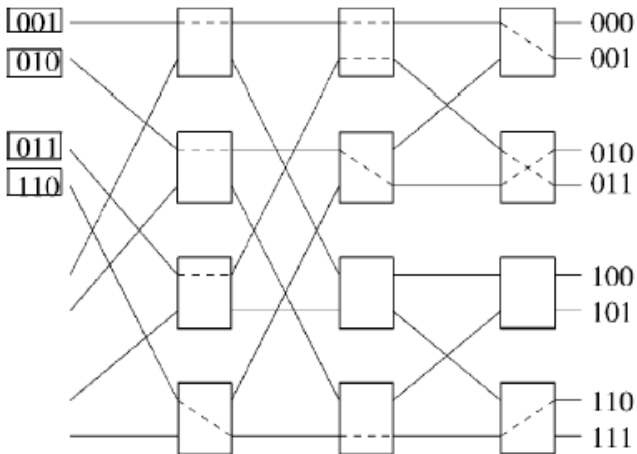
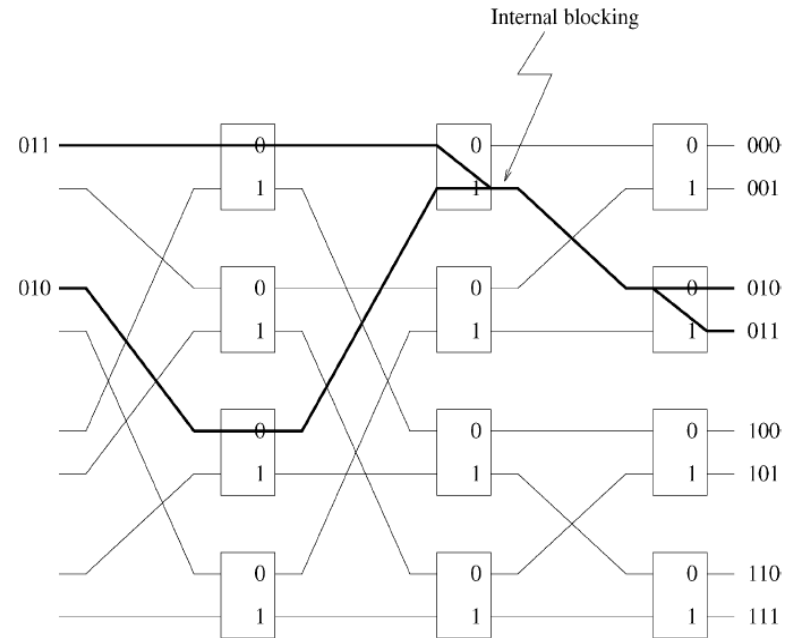
- Stages =  $n = \text{Log}_2N$  ,  
Nodes per stage =  $N/2$
- Self routing: one bit check in each step
- Regularity: Attractive for VLSI implementation



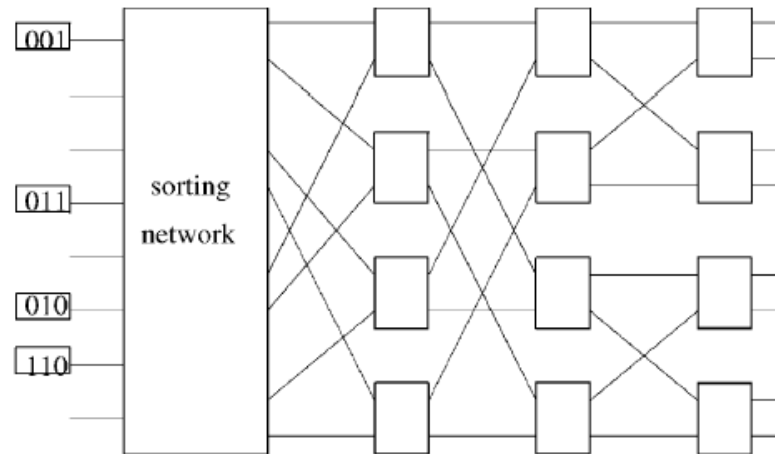
# INTERNAL BLOCKING

## Internal blocking

- Cell is lost due to the contention on a link inside the network
- Can not occur in banyan networks if
  - There is no idle input between any two active inputs
  - Destination address of cells are sorted
- The need for sorting network



(a)

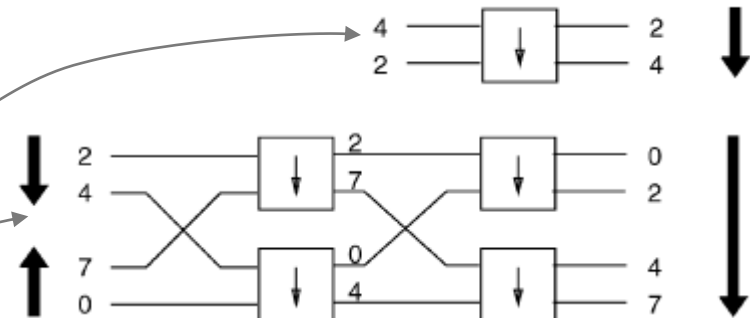


(b)

# BATCHER SORTING NETWORK

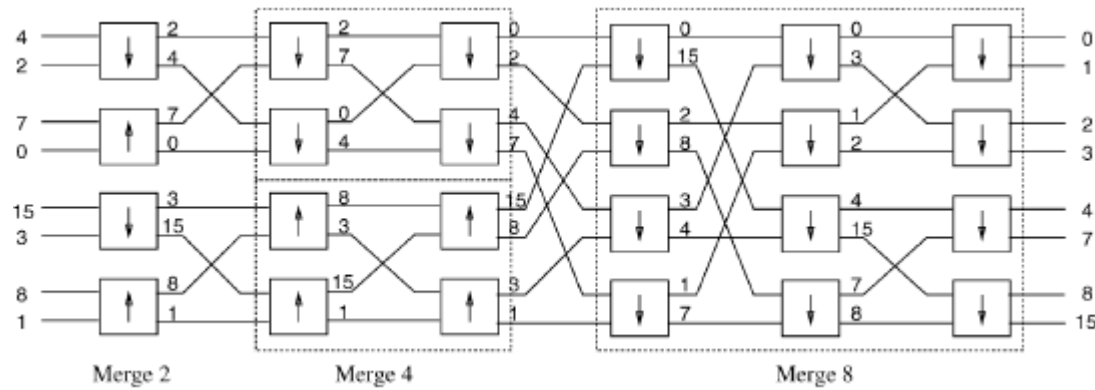
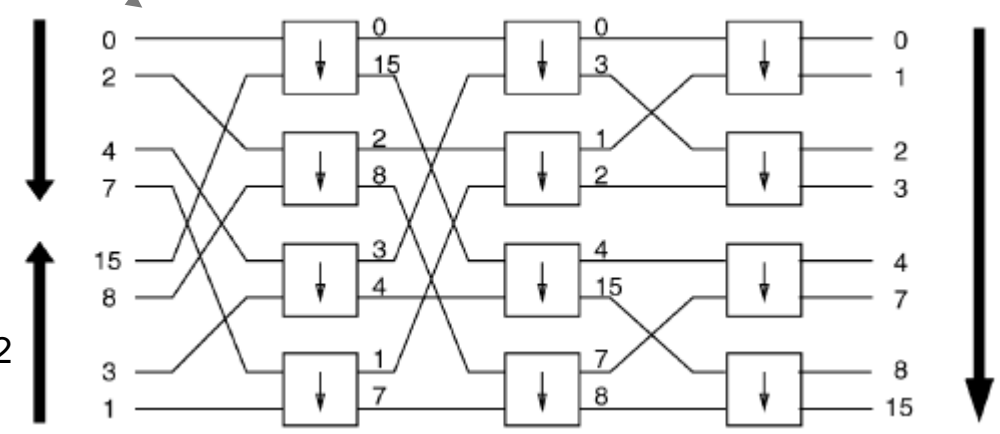
## Merge network

- Consists of  $2 \times 2$  sorting elements
- Partial sort
- Merge2
- Merge4
- Merge8
- MergeN contains  $(N \log_2 N) / 2$  SEs



## Batcher sorting network

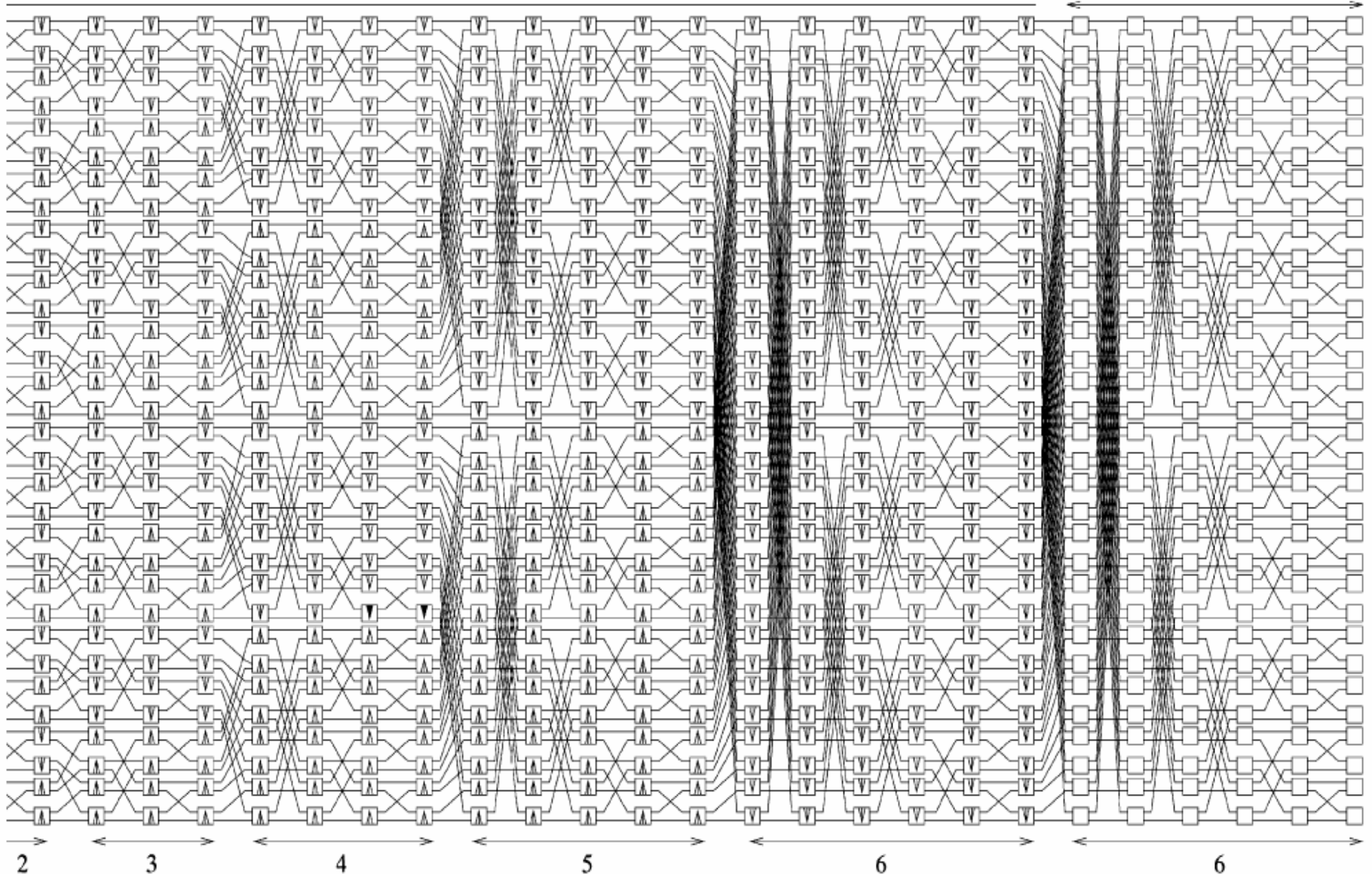
- A series of merge networks
- $8 \times 8$ : merge2  $\rightarrow$  merge4  $\rightarrow$  merge8
- $N \times N$  batcher network
  - $(N \log_2 N)(\log_2 N + 1) / 4$  SEs
  - $1 + 2 + \dots + \log_2 N = (\log_2 N)(\log_2 N + 1) / 2$  Stages



# BATCHER SORTING NETWORK

64 × 64 Batcher network

64 × 64 Banyan network



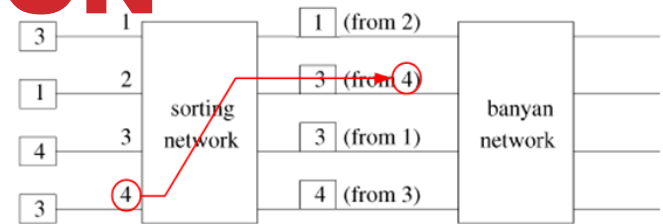
# OUTPUT CONTENTION RESOLUTION

## Output contention resolution algorithms

- Three phase method
- Ring reservation method

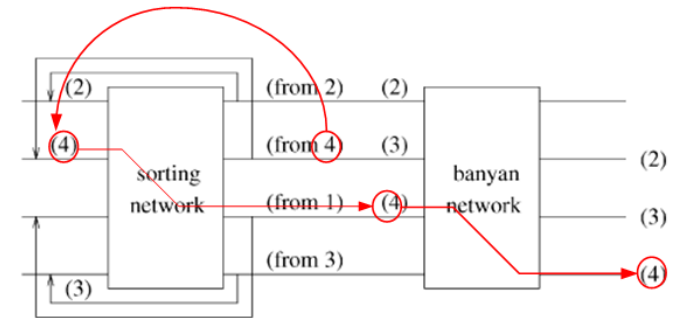
## Three phase method

- phases
  - Request
  - Acknowledge
  - Data



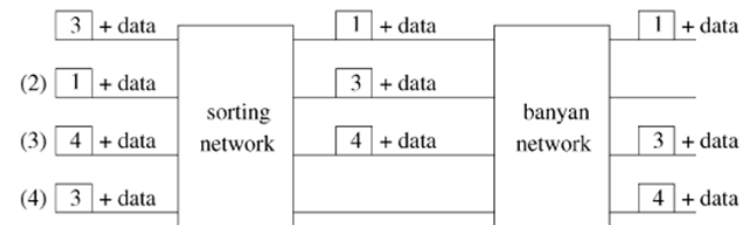
(a) Phase I: Send and resolve request

- Send source-destination pair through sorting network
- Sort destination in non-decreasing order
- Purge adjacent requests with same destination



(b) Phase II: Acknowledge winning ports

- Send ACK with source to ports winning contention
- Route ACK through Batcher-banyan network



(c) Phase III: Send with data

- Acknowledged ports send cells through Batcher-banyan network
- Cells not acknowledged are buffered and retry in the next slot

# OUTPUT CONTENTION RESOLUTION

## Ring reservation method

- A batcher-banyan fabric
- A ring head end (RHE)
- N switch interfaces

## Each switch interface

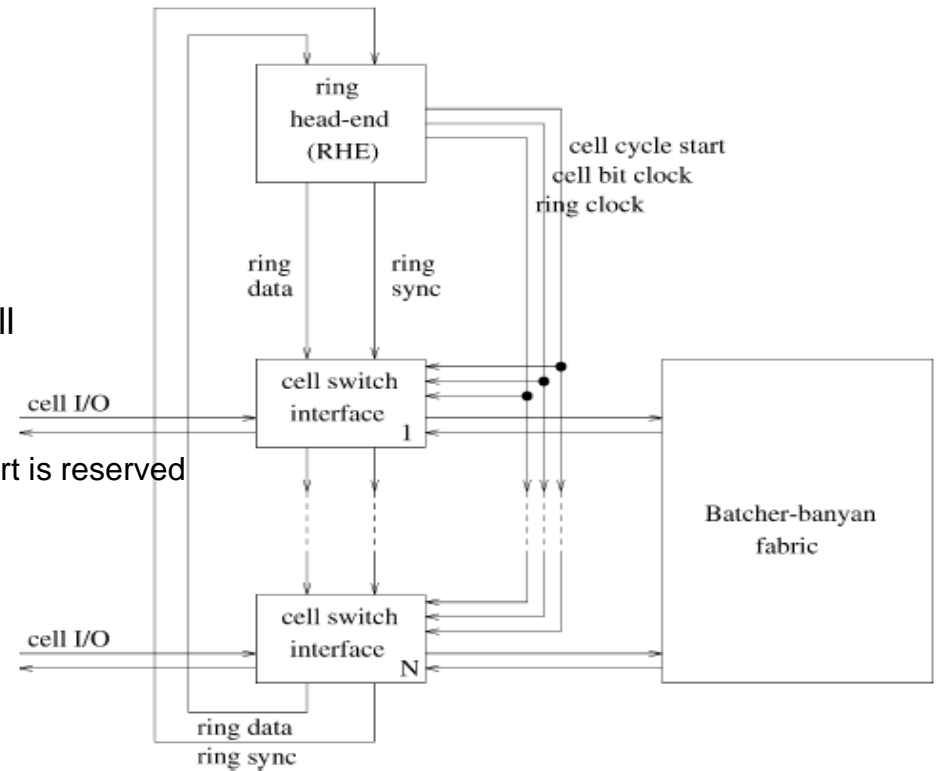
- Buffering input and output cells
- Competition for reserving output port of the cell
- A counter in each switch interface
  - Reset by RHE control signals
  - Incremented at each time slot
  - Match with output port number → output port is reserved

## Advantages

- Fairness
- Can be used for all switch types

## Disadvantage

- N time slots for each arbitration



Input cells	Time slot #1		Time slot #2		Time slot #3		Time slot #4		Time slot #5		Time slot #6		Output cells
	X	Z	X	Z	X	Z	X	Z	X	Z	X	Z	
#0 [ 3 ]	0	0	1	1	0	2	0 ✓	3	0	4	1	5	[ 3 ]
#1 [ 1 ]	0 ✓	1	0	2	0	3	0	4	1	5	0	0	[ 1 ]
#2 [ 5 ]	0	2	0	3	0	4	1 ×	5	0	0	1	1	
#3 [ 1 ]	0	3	0	4	1	5	0	0	1 ×	1	0	2	
#4 [ 3 ]	0	4	1	5	0	0	1	1	0	2	1 ×	3	
#5 [ 5 ]	0 ✓	5	0	0	1	1	0	2	1	3	0	4	[ 5 ]

✓ : allowed to send

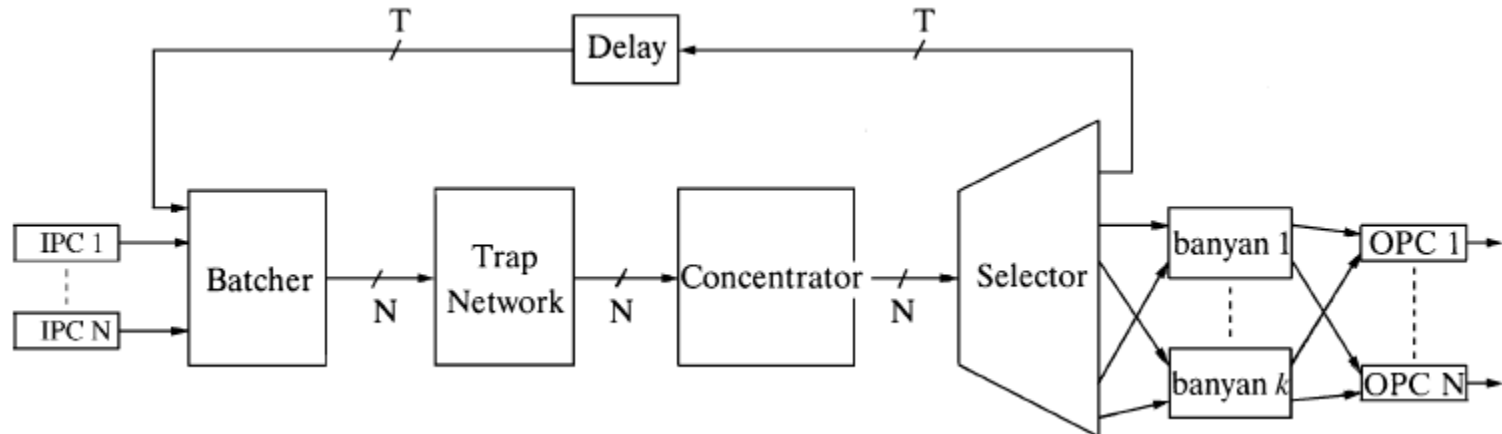
× : not allowed to send



# THE SUNSHINE SWITCH

## The Sunshine switch

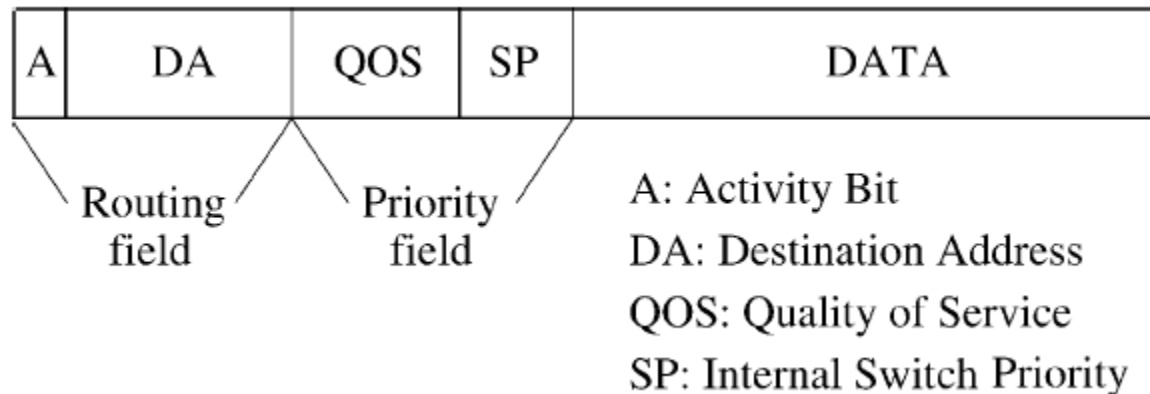
- A batcher sorting network
- $k$  banyan networks in parallel
- More than one path to each destination ( $k$  paths)
- Recirculating queue
  - $T$  paths in the figure
  - If more than  $k$  cells have the same output port
  - Excess cells are recirculated with a delay
- Batcher network
  - Sorts in the order of port & priority
  - The highest priority cell for each port is selected by trap network



# THE SUNSHINE SWITCH

## The sunshine switch

- Control Header
  - Is added to each cell by IPC
  - Two parts
    - Routing part
      - A: cell activity (non-emptiness)
      - DA: destination address
    - Priority part
      - QoS: quality of service (priority of cell)
      - SP: switch priority (assigned by switch to compensate the recirculation delay)



# DEFLECTION ROUTING

## Deflection

- Two cells contend at a node
- One of them will be routed incorrectly

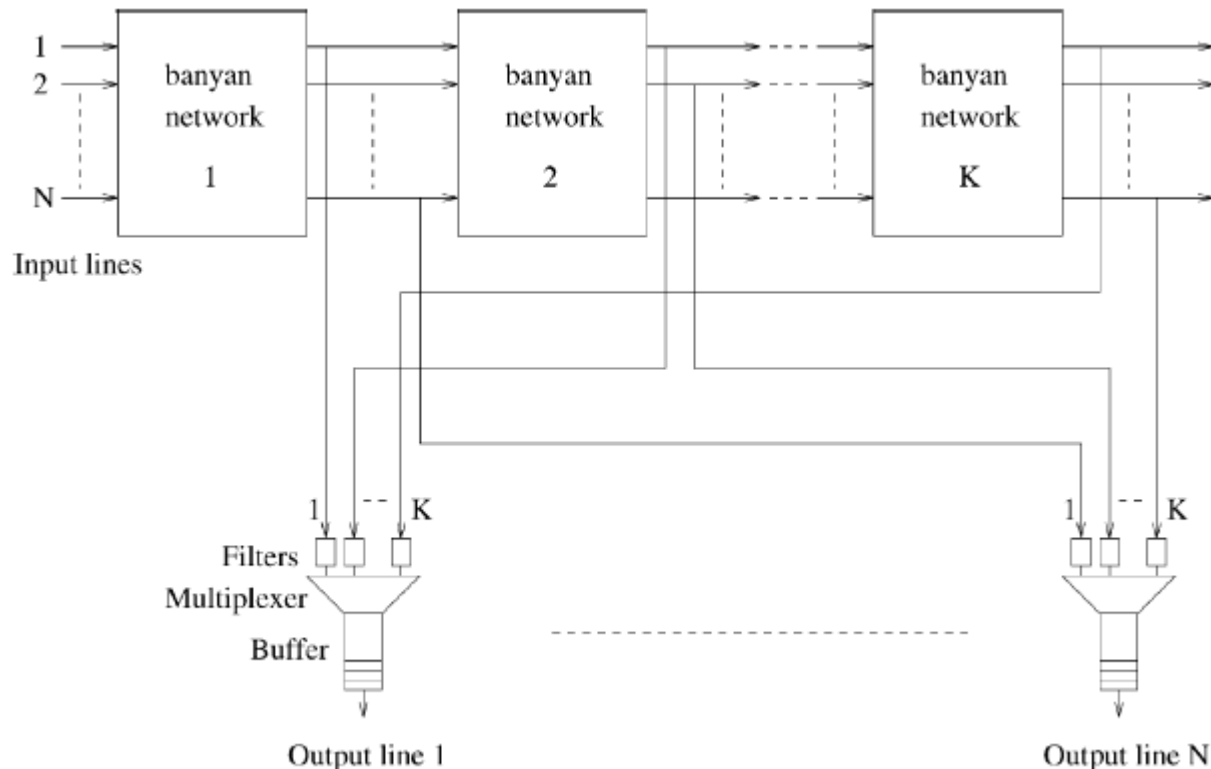
## Works on deflection routing

- Tandem banyan switch
- Shuffle-exchange network with deflection routing
- Dual shuffle exchange network with error-correcting routing

# TANDEM BANYAN SWITCH

## Tandem banyan switch

- Chain of K banyan networks
- Deflected cells continue into the next banyan network
- Correctly routed cells go to output buffers
- One added banyan network  $\rightarrow$  one order of magnitude reduction in deflections



# TANDEM BANYAN SWITCH

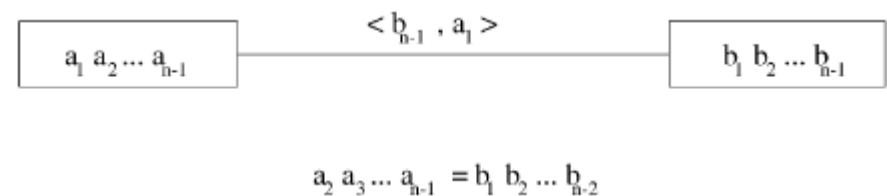
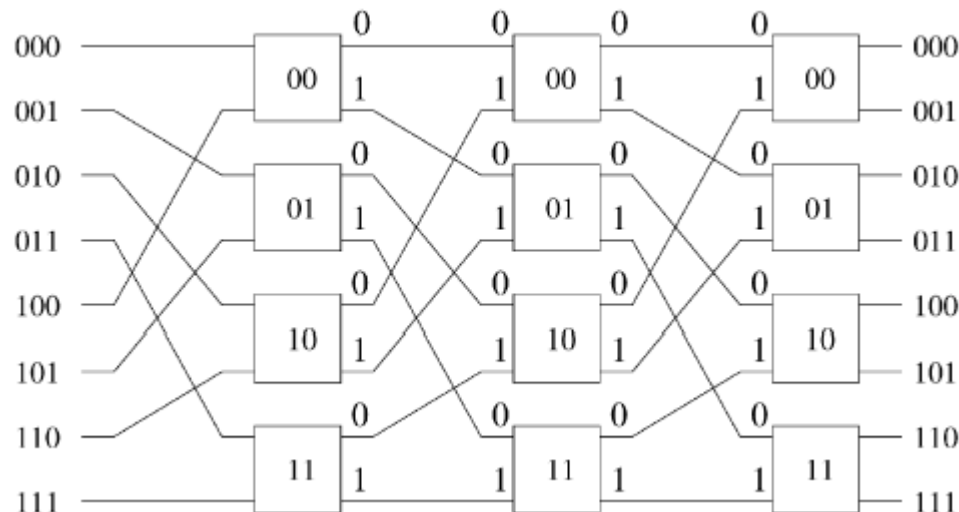
## Tandem banyan switch

- Switching header
  - Activity bit:  $a$
  - Conflict bit:  $c$
  - Priority field:  $P$
  - Address field:  $D$
- For two cells 1, 2 in the same stage  $s$ :
  1. If  $a_1 = a_2 = 0$ , then take no action, i.e., leave the switch in the present state.
  2. If  $a_1 = 1$  and  $a_2 = 0$ , then set the switch according to  $d_{s1}$ .
  3. If  $a_1 = 0$  and  $a_2 = 1$ , then set the switch according to  $d_{s2}$ .
  4. If  $a_1 = a_2 = 1$ , then:
    - (a) If  $c_1 = c_2 = 1$ , then take no action.
    - (b) If  $c_1 = 0$  and  $c_2 = 1$ , then set the switch according to  $d_{s1}$ .
    - (c) If  $c_1 = 1$  and  $c_2 = 0$ , then set the switch according to  $d_{s2}$ .
    - (d) If  $c_1 = c_2 = 0$ , then:
      - i. If  $P_1 > P_2$ , then set the switch according to  $d_{s1}$ .
      - ii. If  $P_1 < P_2$ , then set the switch according to  $d_{s2}$ .
      - iii. If  $P_1 = P_2$ , then set the switch according to either  $d_{s1}$  or  $d_{s2}$ .
      - iv. If one of the cells has been misrouted, then set its conflict bit to 1

# SHUFFLE-EXCHANGE NETWORK WITH DEFLECTION ROUTING

## N\*N shuffle-exchange network (SN)

- Stages =  $n = \log_2 N$
- SEs per stage =  $N/2$
- SE labels: numbers with  $n-1$  bits length
- SE input (output) labels: 1 bit numbers
- How SE forwards cells
  - The cell with a 0 in  $i$ 'th destination address bit goes to output 0
  - The cell with a 1 in  $i$ 'th destination address bit goes to output 1
- Network connections:
  - Consider binary labels represented in form  $(a_1 a_2 \dots)$
  - Output  $a_n$  of node  $X=(a_1 a_2 \dots a_{n-1}) \rightarrow$  input  $a_1$  of node  $Y=(a_2 a_3 \dots a_n)$
  - This link is represented as  $\langle a_n, a_1 \rangle$



# SHUFFLE EXCHANGE NETWORK WITH DEFLECTION ROUTING

## The path from input to output

- Is determined by:
  - Source address  $s_1 s_2 \dots s_n$
  - Destination address  $d_1 d_2 \dots d_n$

$$\begin{array}{rcll}
 S = s_1 \dots s_n & & & \\
 \xrightarrow{\langle -, s_1 \rangle} & (s_2 \dots s_n) & \xrightarrow{\langle d_1, s_2 \rangle} & (s_3 \dots s_n d_1) \\
 \xrightarrow{\langle d_2, s_3 \rangle} & \dots & \xrightarrow{\langle d_{i-1}, s_i \rangle} & (s_{i+1} \dots s_n d_1 \dots d_{i-1}) \\
 \xrightarrow{\langle d_i, s_{i+1} \rangle} & \dots & \xrightarrow{\langle d_{n-1}, s_n \rangle} & (d_1 \dots d_{n-1}) \\
 \xrightarrow{\langle d_n, 0 \rangle} & d_1 \dots d_n = D. & & 
 \end{array}$$

- An example: S=001 and D=101

$$\begin{array}{ccccccccc}
 001 & \xrightarrow{\langle -, 0 \rangle} & 01 & \xrightarrow{\langle 1, 0 \rangle} & 11 & \xrightarrow{\langle 0, 1 \rangle} & 10 & \xrightarrow{\langle 1, 0 \rangle} & 101 \\
 S & \xrightarrow{\langle -, s_1 \rangle} & s_2 s_3 & \xrightarrow{\langle d_1, s_2 \rangle} & s_3 d_1 & \xrightarrow{\langle d_2, s_3 \rangle} & d_1 d_2 & \xrightarrow{\langle d_3, 0 \rangle} & D
 \end{array}$$

- SEs the cell passes
  - An string  $s_2 \dots s_n d_1 \dots d_{n-1}$
  - An  $(n-1)$  bits window shifting one bit in each stage from left to right

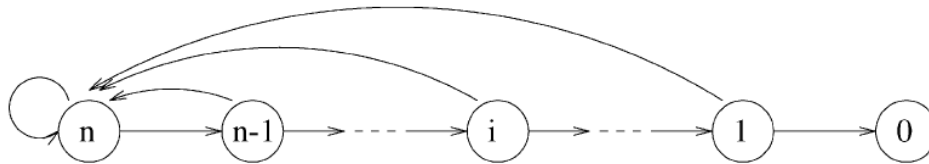




# DUAL SHUFFLE-EXCHANGE NETWORK WITH ERROR-CORRECTING ROUTING

## SN with deflection routing (the previous scheme)

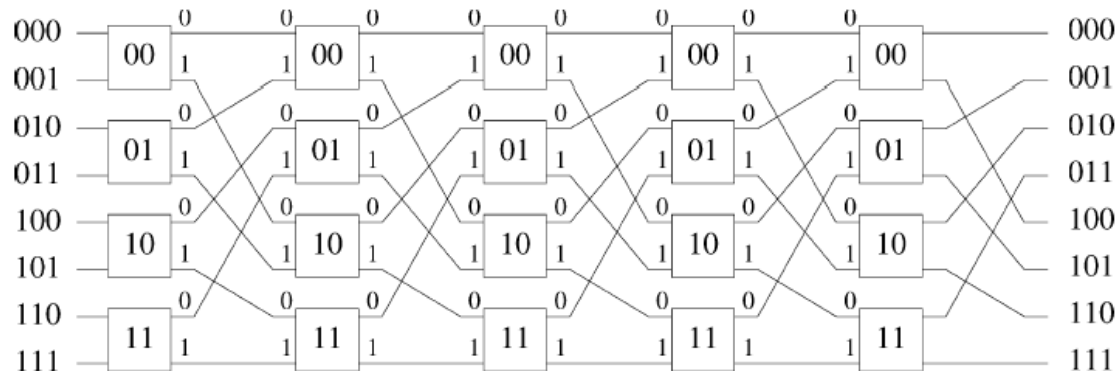
- Highly inefficient: Routing must be restarted for deflected cell:



- Desired network behavior:



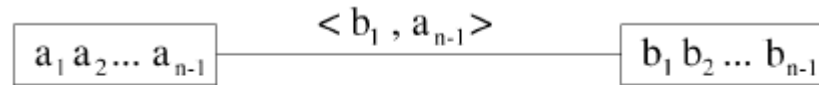
- Dual shuffle-exchange network
  - A shuffle-exchange network (SN)
  - An unshuffle-exchange network (USN)
    - Mirror of SN network
    - An 8\*8 example of USN:



# DUAL SHUFFLE-EXCHANGE NETWORK WITH ERROR-CORRECTING ROUTING

USN: mirror image of SN

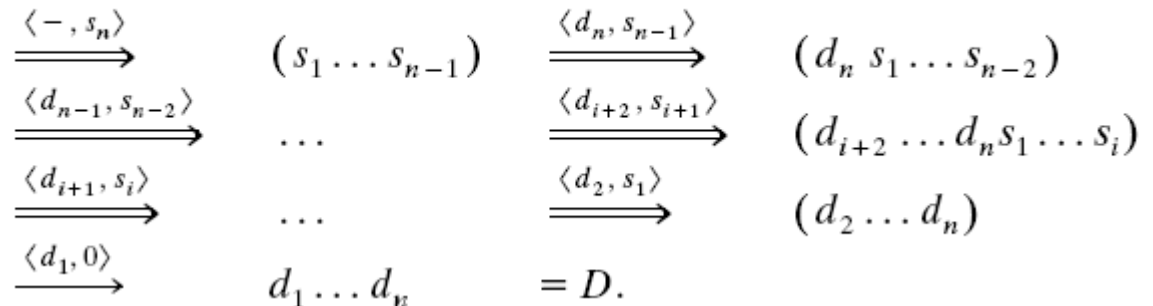
Similar rules as SN:



$$a_1 a_2 \dots a_{n-2} = b_2 b_3 \dots b_{n-1}$$

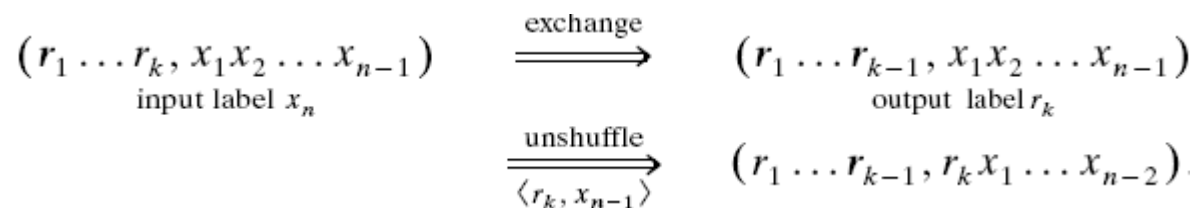
- SE connections:
- Paths from inputs to outputs:

$$S = s_1 \dots s_n$$



- SEs the cell passes
  - An string  $d_2 \dots d_n s_1 \dots s_{n-1}$
  - An  $(n-1)$  bits window shifting one bit in each stage from right to left
- Traveling cell state diagram:

- Initial state:  $(d_1 \dots d_n, s_1 \dots s_{n-1})$
- Final state:  $(d_1 \dots d_n)$
- Transition:



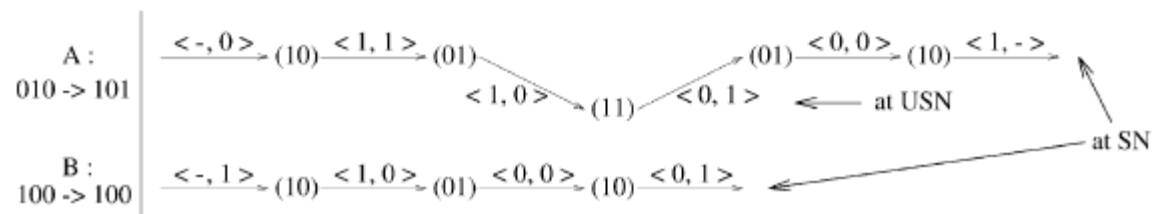
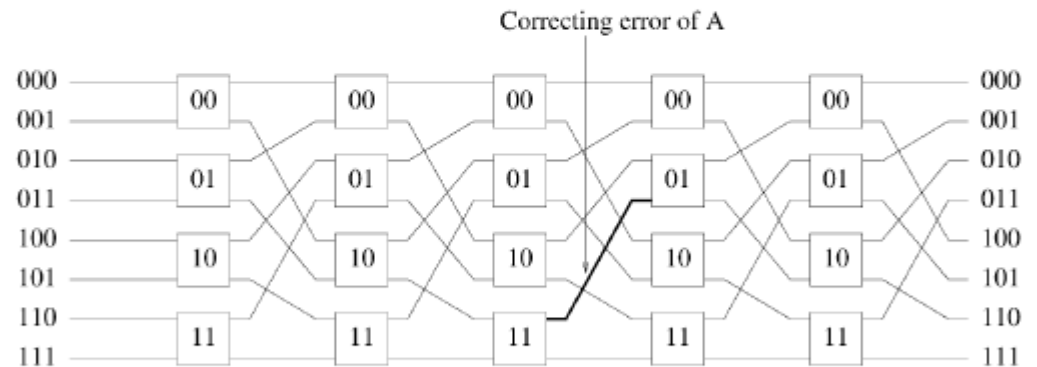
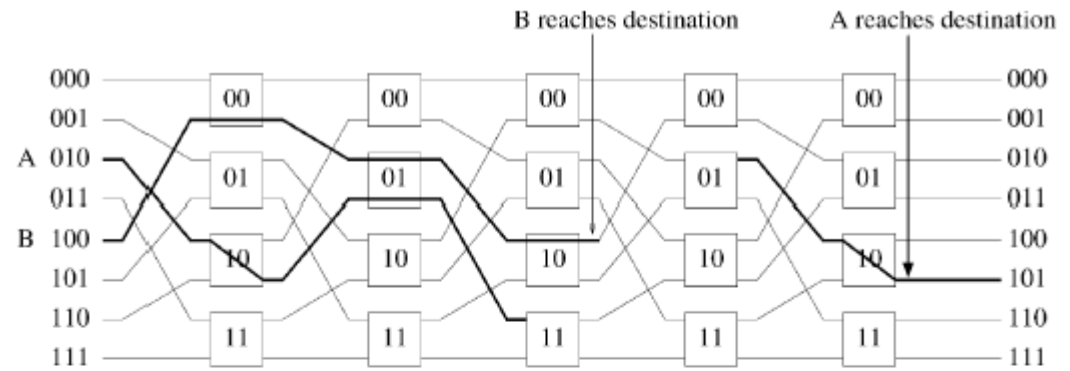
# DUAL SHUFFLE-EXCHANGE NETWORK WITH ERROR-CORRECTING ROUTING

Consider a USN overlaid on a SN

USN can undo what SN performs

Deflected cell can return to the state before deflection

Example:

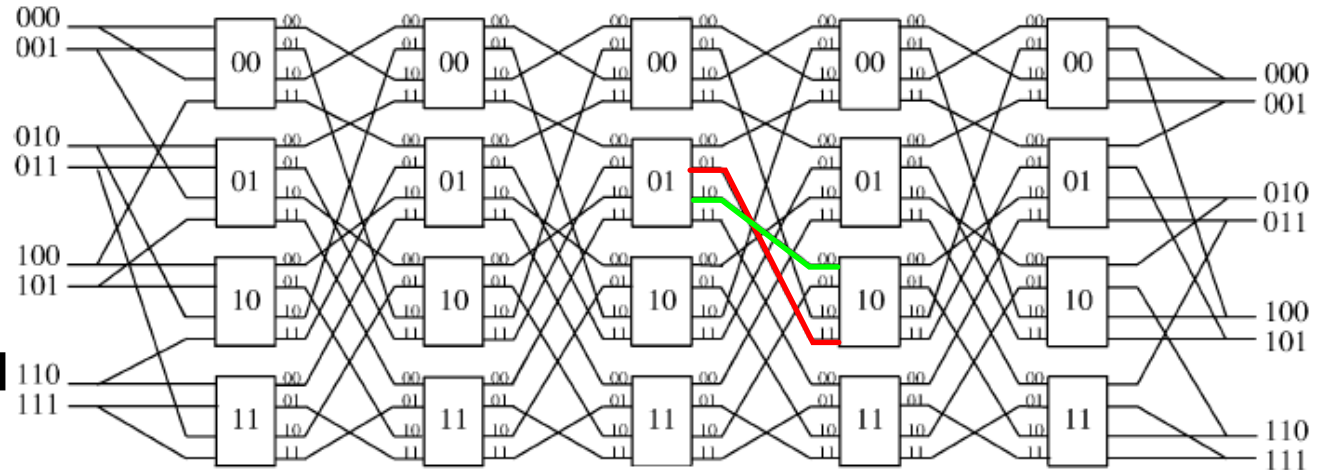


# DUAL SHUFFLE-EXCHANGE NETWORK WITH ERROR-CORRECTING ROUTING

## Error correction procedure

- Cell state:  $(r_1 \dots r_k, x_1 \dots x_{n-1})$
- Cell should be sent via link  $\langle r_k, x_1 \rangle$  to the next stage
- It is deflected to link  $\langle r'_k, x_1 \rangle$
- It goes to node  $(x_2 \dots x_{n-1} r'_k)$  in the next stage
- Error correction procedure does not remove the bit  $r_k$
- Instead, it attaches bit  $x_1$  to the routing tag
- New state:  $(r_1 \dots r_k x_1, x_2 \dots x_{n-1} r'_k)$
- Then the cell is moved to the USN
- It will be sent via link  $\langle x_1, r'_k \rangle$
- It will return to the state before deflection:  $(r_1 \dots r_k, x_1 \dots x_{n-1})$
- If cell is deflected in the USN, it can be corrected in SN in a similar way

# DUAL SHUFFLE-EXCHANGE NETWORK WITH ERROR-CORRECTING ROUTING



## Merging SN and USN

- SN: 2\*2 SEs
- USN: 2\*2 SEs
- Merged (dual SN): 4\*4 SEs

## Labeling

- Inputs: 00..11
- Outputs: 00..11
- Unshuffle links:  $\langle 0a, 1b \rangle$ 
  - Connect outputs 10 or 11 to inputs 00 or 01 of the next stage
- Shuffle links:  $\langle 1a, 0b \rangle$ 
  - Connect outputs 00 or 01 to inputs 10 or 11 of the next stage

Nodes interconnected by a shuffle link

$$a_1 a_2 \dots a_{n-1} \xrightarrow{\langle 1b_{n-1}, 0a_1 \rangle} b_1 b_2 \dots b_{n-1} \quad a_2 \dots a_{n-1} = b_1 \dots b_{n-1}$$

Nodes interconnected by an unshuffle link

$$a_1 a_2 \dots a_{n-1} \xrightarrow{\langle 0b_1, 1a_{n-1} \rangle} b_1 b_2 \dots b_{n-1} \quad a_1 \dots a_{n-2} = b_2 \dots b_{n-2}$$

## Connections

- Consider two nodes  $A=(a_1 \dots a_{n-1})$ ,  $B=(b_1 \dots b_{n-1})$
- They are connected via unshuffle link  $\langle 0b_1, 1a_{n-1} \rangle$  if  $a_1 \dots a_{n-2} = b_2 \dots b_{n-1}$
- They are connected via shuffle link  $\langle 1b_{n-1}, 0a_1 \rangle$  if  $a_2 \dots a_{n-1} = b_1 \dots b_{n-2}$

# DUAL SHUFFLE-EXCHANGE NETWORK WITH ERROR-CORRECTING ROUTING

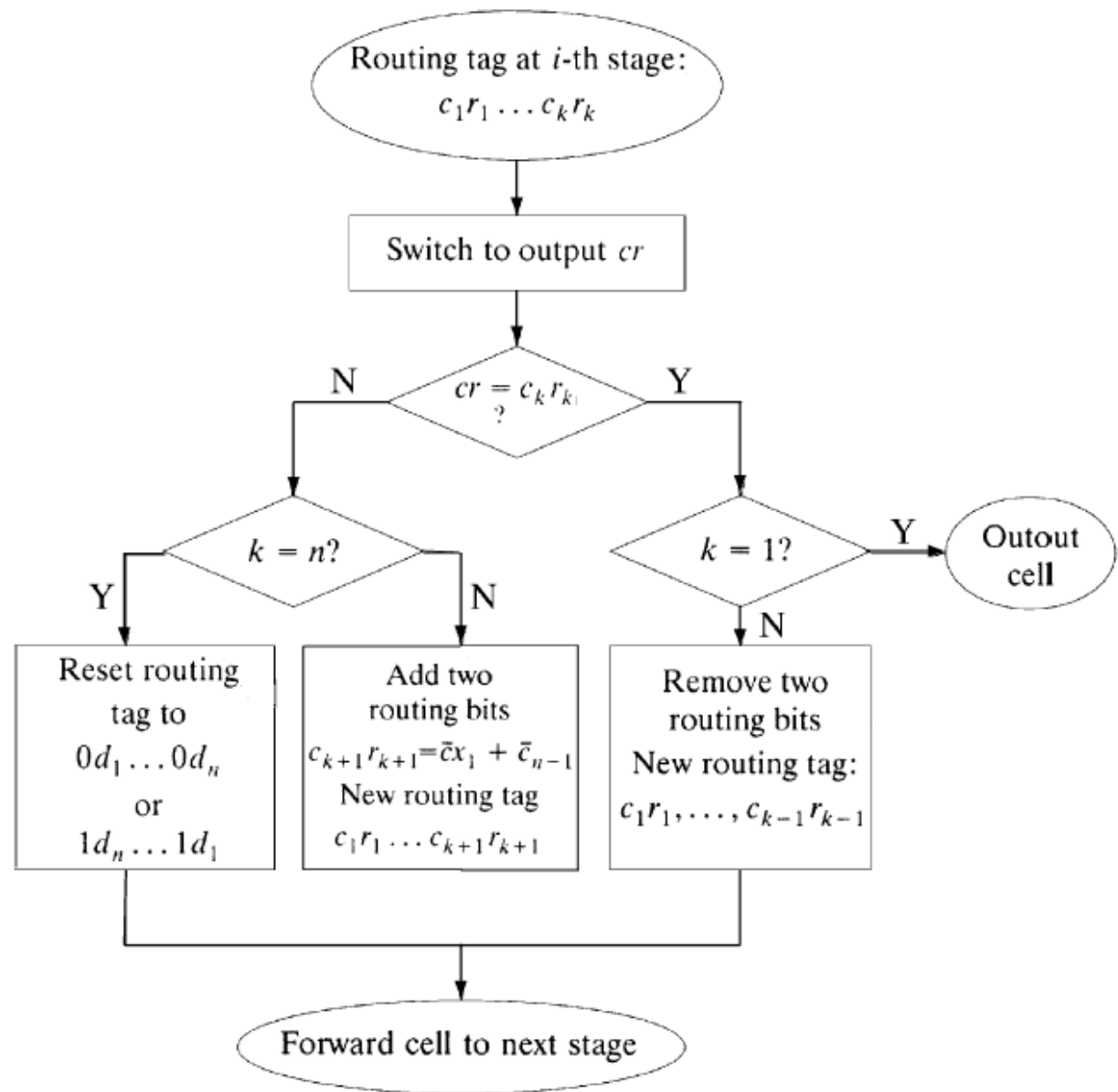
## States and transitions

- Two routing bits required at each stage
- Destination:  $D=d_1\dots d_n$
- Initial state:
  - $0d_1\dots 0d_n$  (Starting from USN)
  - $1d_n\dots 1d_1$  (Starting from SN)
- Cell state:  $(c_1r_1\dots c_kr_k, x_1\dots x_{n-1})$
- Transitions:
  - $c_k=0 \rightarrow$  unshuffle link
  - $c_k=1 \rightarrow$  shuffle link
  -

$$(c_1r_1 \dots c_kr_k, x_1 \dots x_{n-1})$$
$$\left\{ \begin{array}{ll} \xrightarrow{\langle 0r_k, 1x_{n-1} \rangle} & (c_1r_1 \dots c_{k-1}r_{k-1}, r_k x_1 \dots x_{n-2}) \quad \text{if } c_k = 0, \\ \xrightarrow{\langle 1r_k, 0x_1 \rangle} & (c_1r_1 \dots c_{k-1}r_{k-1}, x_2 \dots x_{n-1}r_k) \quad \text{if } c_k = 1. \end{array} \right.$$

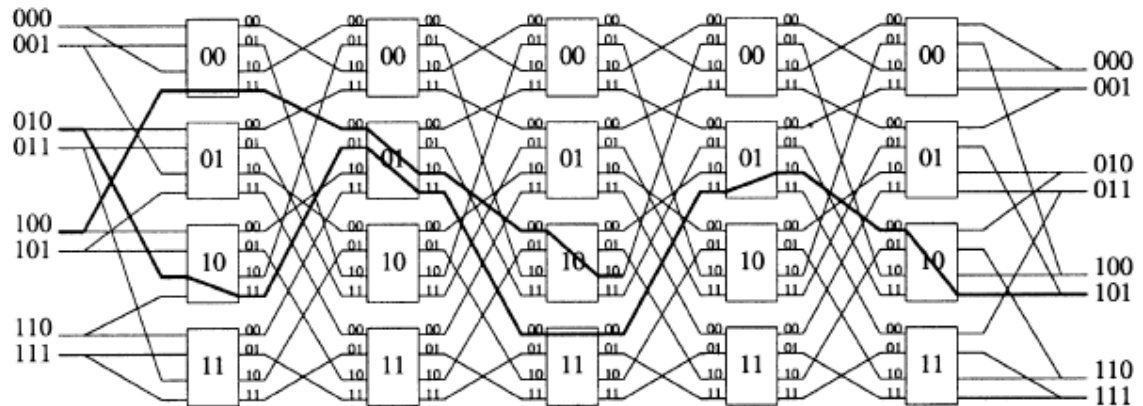
# DUAL SHUFFLE-EXCHANGE NETWORK WITH ERROR-CORRECTING ROUTING

## The algorithm



# DUAL SHUFFLE-EXCHANGE NETWORK WITH ERROR-CORRECTING ROUTING

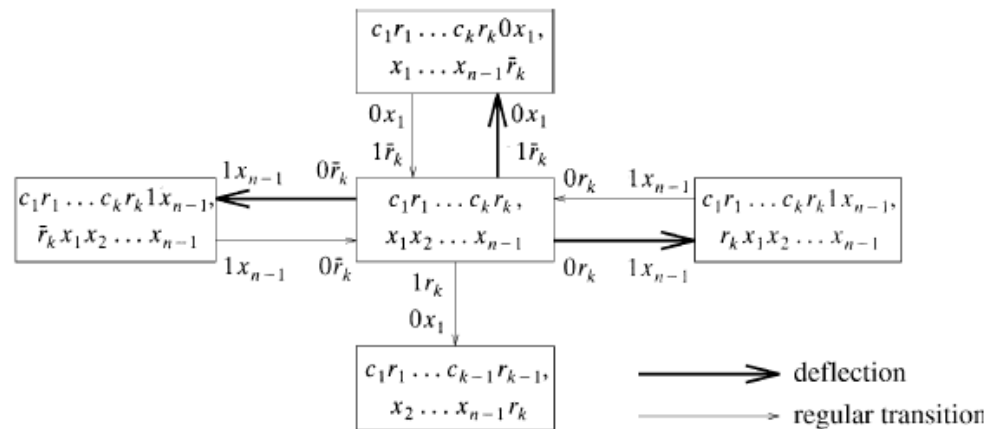
An example  
(The previous example)



State transition of cells A and B

State diagram

error      error correction  
 A: 010 → (111011, 10) → (1110, 01) → (111000, 11) → (1110, 01) → (11, 10) → 101  
 B: 100 → (101011, 00) → (1010, 01) → (10, 10) → 100



Transitions  
for deflection

$$\begin{aligned}
 & (c_1 r_1 \dots c_k r_k, x_1 \dots x_{n-1}) \\
 & \left\{ \begin{array}{l} \xrightarrow{\langle 0r, 1x_{n-1} \rangle} (c_1 r_1 \dots c_k r_k 1x_{n-1}, r x_1 \dots x_{n-2}) \quad \text{if } c_k r_k \neq 0r, \\ \xrightarrow{\langle 1r, 0x_1 \rangle} (c_1 r_1 \dots c_k r_k 0x_1, x_2 \dots x_{n-1} r) \quad \text{if } c_k r_k \neq 1r. \end{array} \right.
 \end{aligned}$$



# MULTICAST COPY NETWORKS

## Point-to-multipoint communications

- A copy network (replicates cells)
- A point-to-point switch

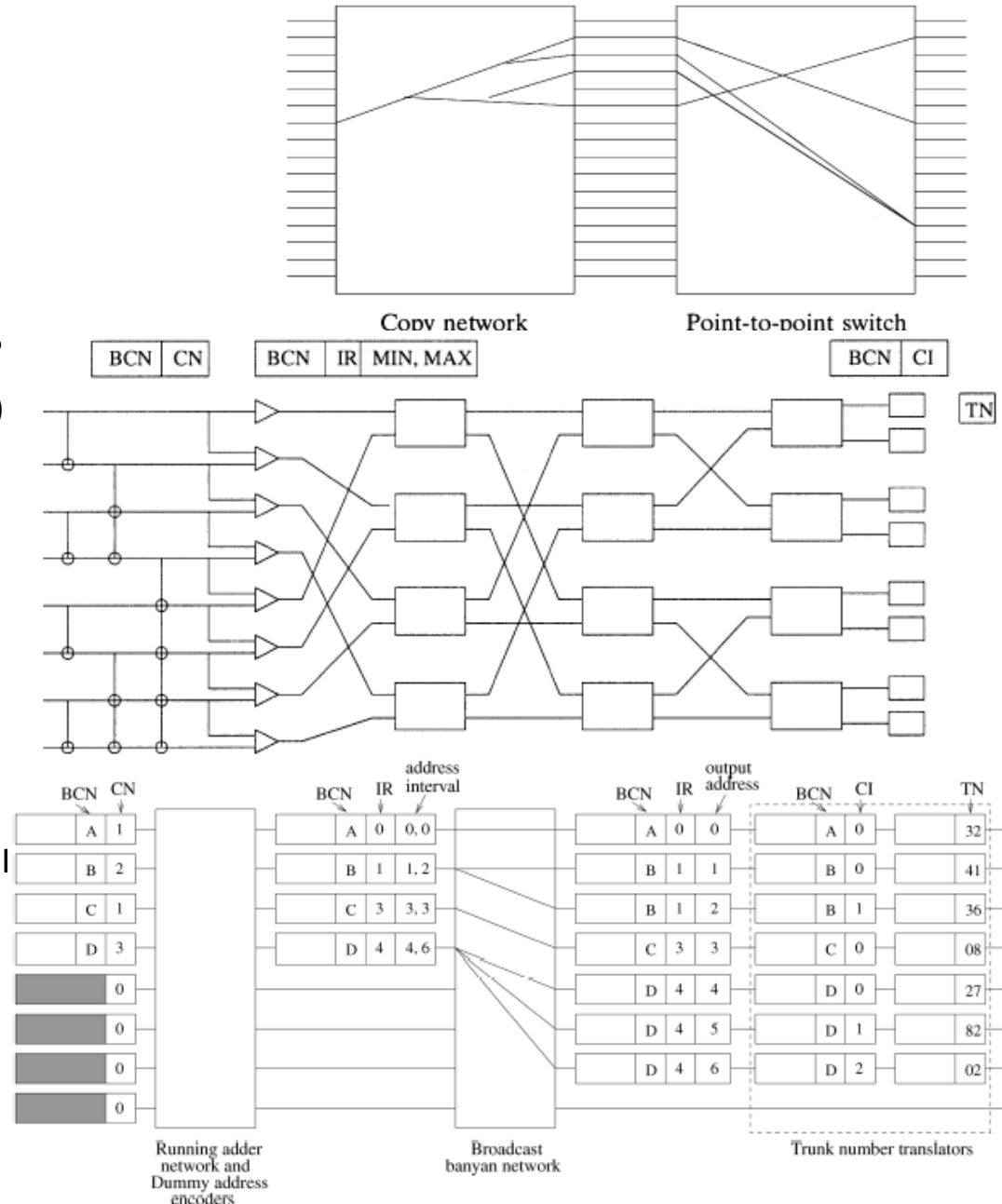
## The copy network components

- Running address network (RAN)
- Dummy address encoder (DAE)
- Broadcast banyan network
- Trunk number translator

## Acronyms in figure:

- CN: Number of copies
- IR: Index reference
- CI: Copy index
- BCN: Broadcast channel number

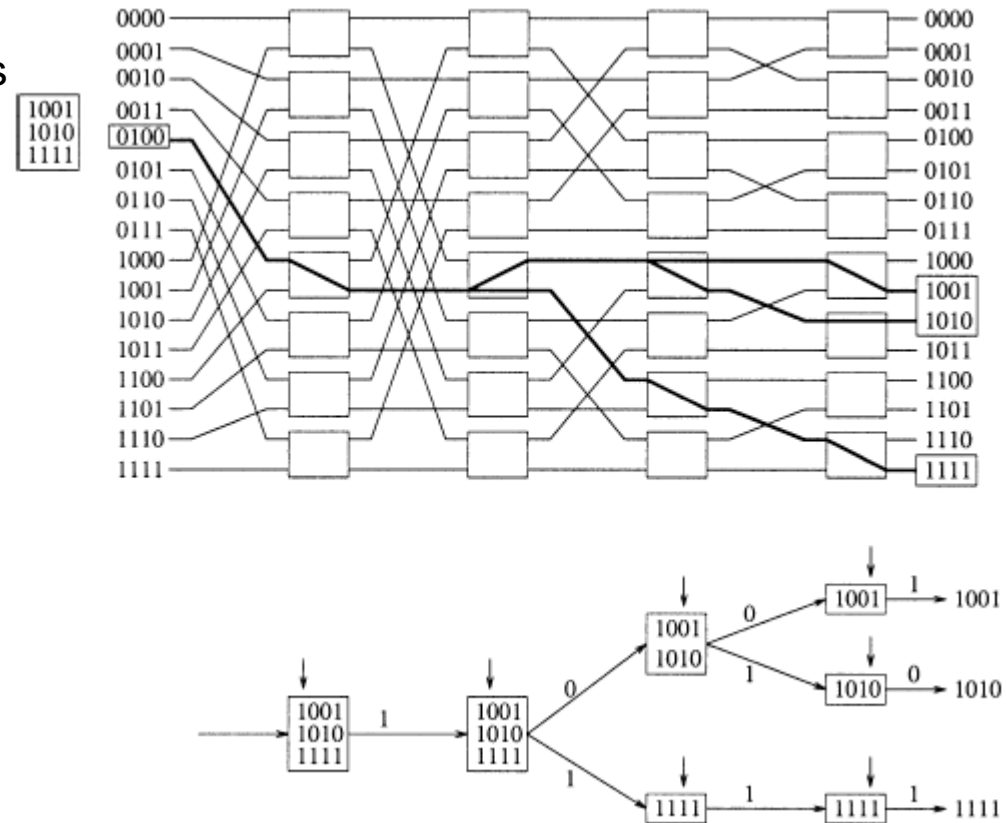
## An example



# MULTICAST COPY NETWORKS

## Broadcast banyan network

- SEs can replicate cells
- 3 possibilities (2 bits needed in cell header):
  - Cell goes to output 1
  - Cell goes to output 1
  - Cell is replicated to both outputs
- Generalized self routing algorithm
  - Current bit of all destination addresses are checked at each stage
  - All zero → cell goes through output 0
  - All one → cell goes through output 1
  - Some zero, some one → cell is replicated
  - An example



# MULTICAST COPY NETWORKS

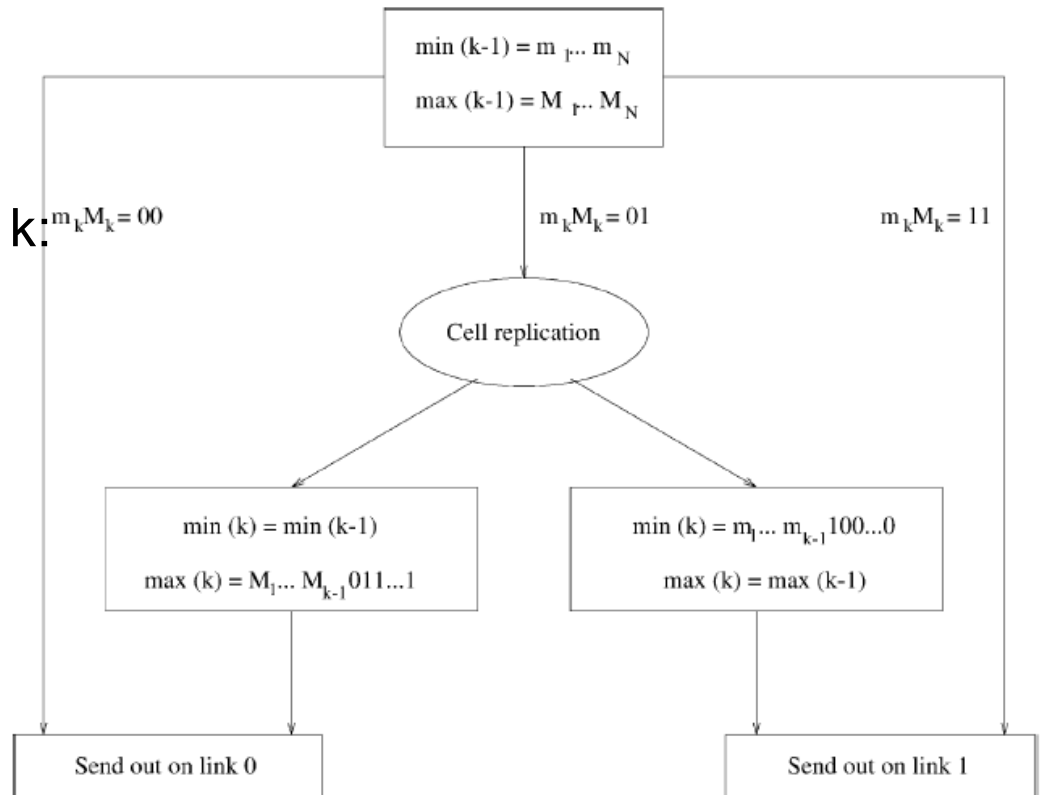
## Broadcast banyan network

- Generalized self routing algorithm
  - Problems
    - A variable number of destination addresses
      - Must be recorded in cell header
      - Must be checked at SEs at each stage
      - Need to be processed for cell header modifications
    - Cell path forms a tree → more blocking
  - Solution
    - Boolean interval splitting algorithm

# MULTICAST COPY NETWORKS

## Boolean interval splitting algorithm

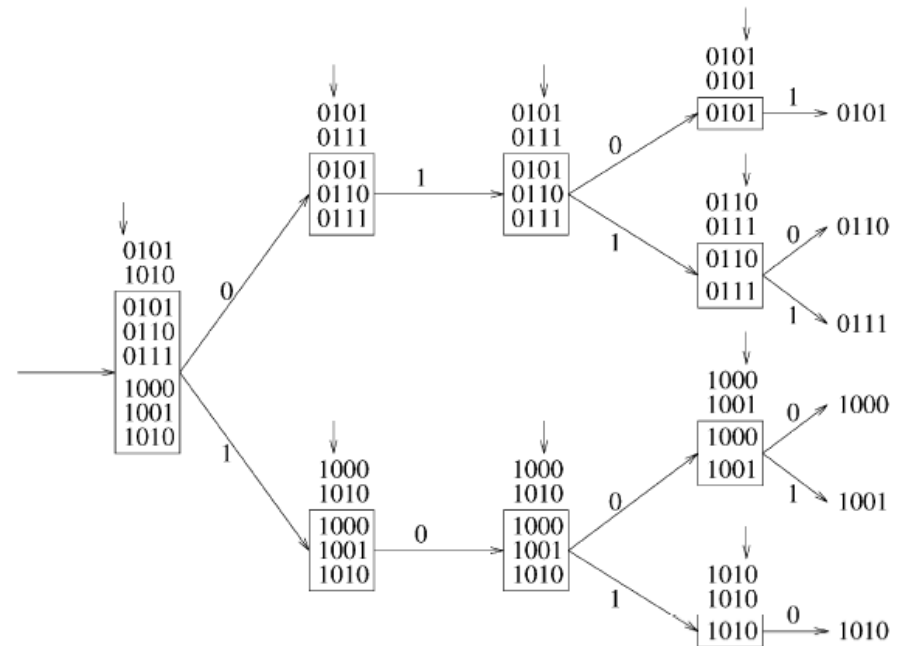
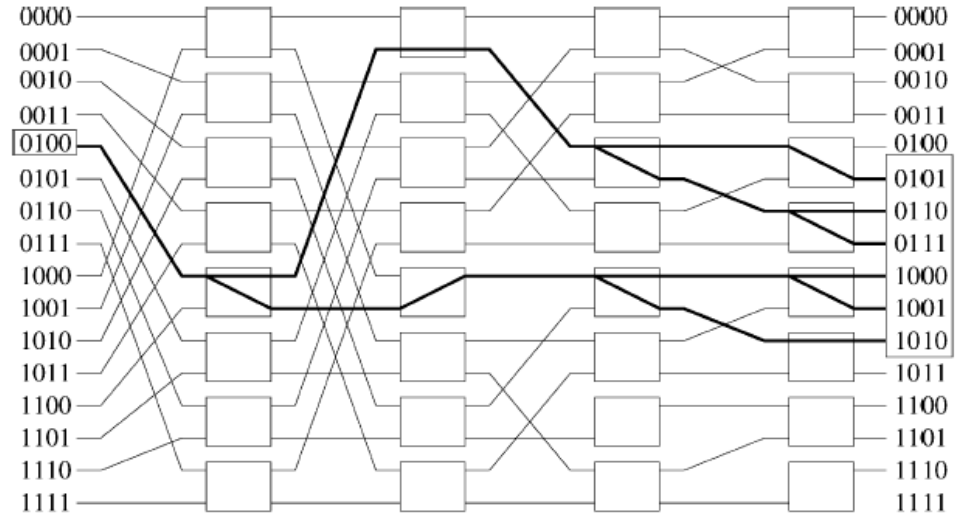
- Interval: a set of N-bits numbers between min and max
- Interval represents a set of contiguous destination addresses
- Stage k
  - $\min(k-1) = m_1 \dots m_k$
  - $\max(k-1) = M_1 \dots M_k$
  - Self routing at stage k:



# MULTICAST COPY NETWORKS

## Boolean interval splitting algorithm

- An example

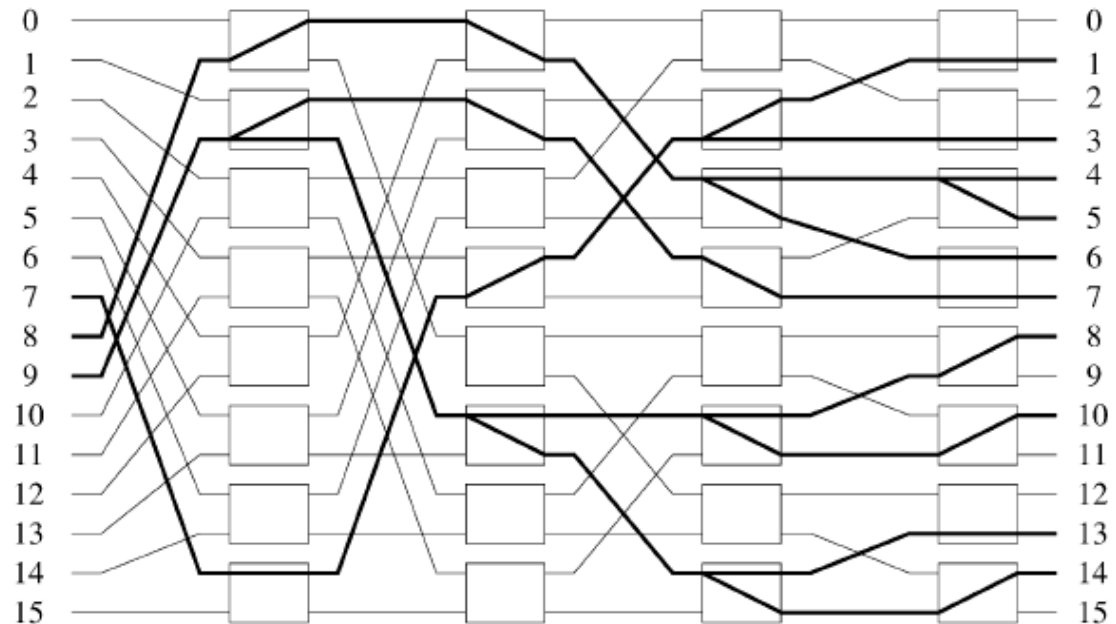


# MULTICAST COPY NETWORKS

Broadcast banyan network is nonblocking when input cells are

- Monotonic
  - Output port sets are sorted
- Concentrated
  - No idle inputs exists between active outputs
- An example

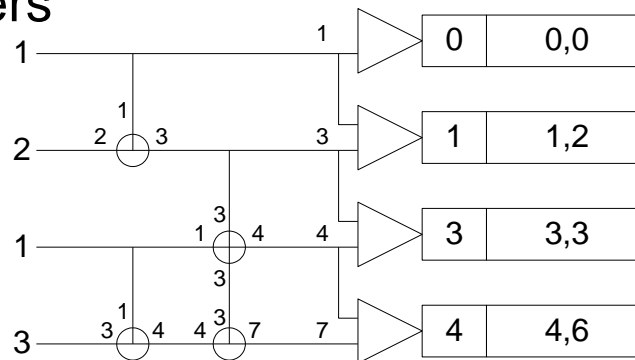
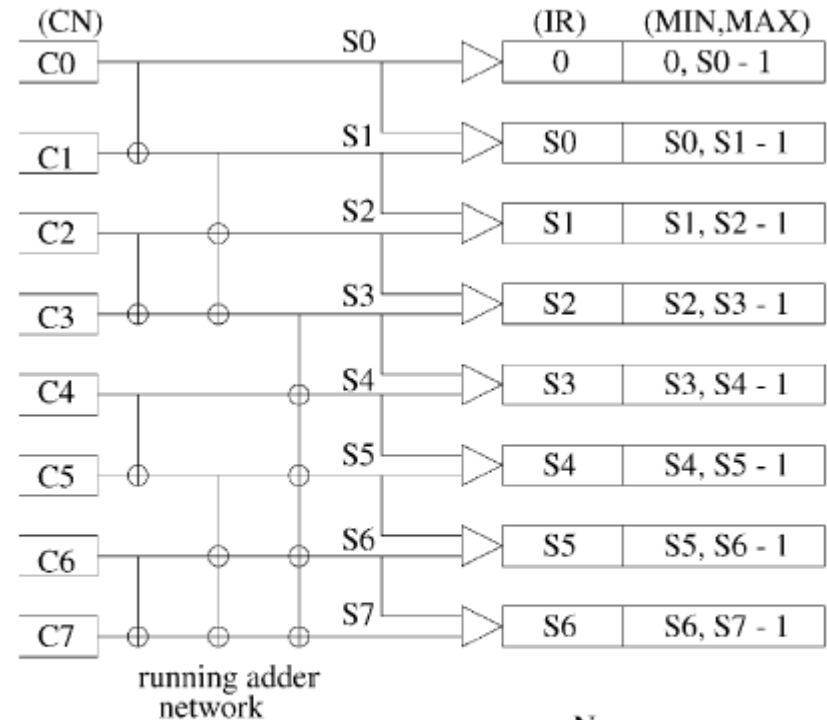
$$\begin{aligned}x_1 &= 7 & Y_1 &= \{1,3\} \\x_2 &= 8 & Y_2 &= \{4,5,6\} \\x_3 &= 9 & Y_3 &= \{7,8,10,13,14\}\end{aligned}$$



# MULTICAST COPY NETWORKS

## RAN, DAE

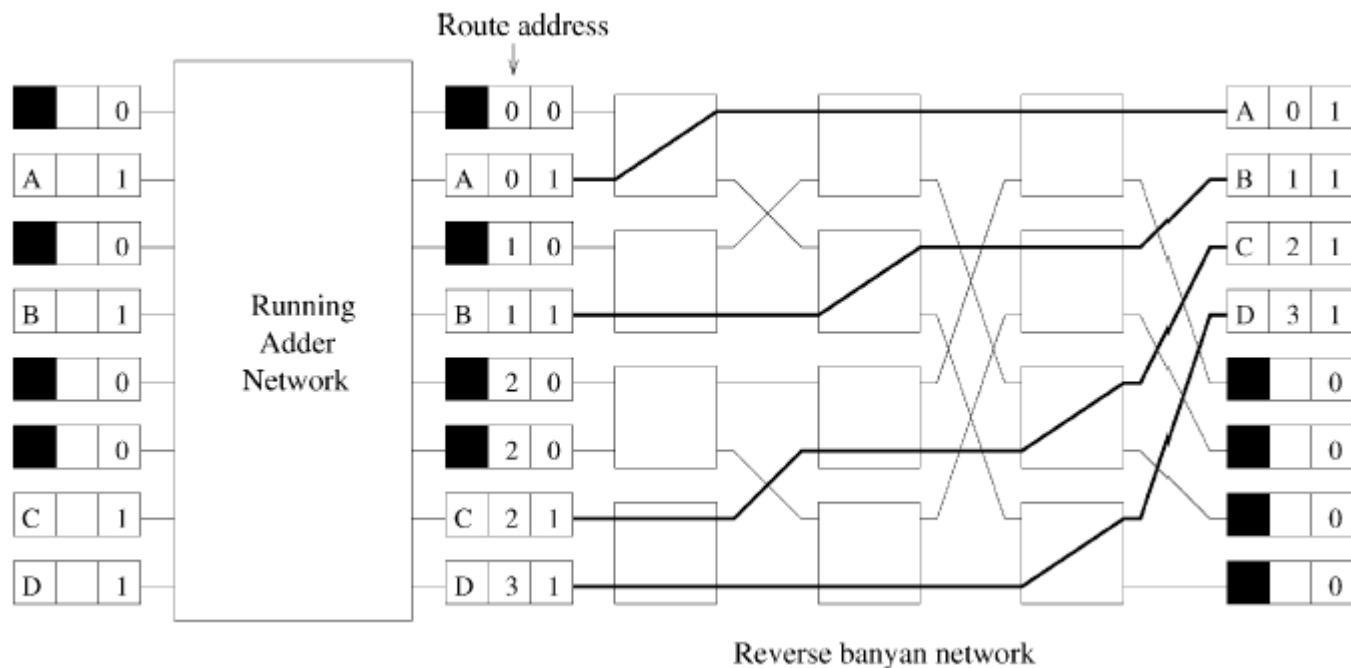
- Encoding process
  - RAN, DAE transform copy numbers into a set of monotonic addresses
- Decoding process
  - TNT determines the real destinations of copies
- RAN, DAE
  - $(N/2)\log_2 N$  adders
  - $\log_2 N$  stages
- An example



# MULTICAST COPY NETWORKS

## Concentration

- Eliminate idle inputs between active outputs
- Must be done before broadcast banyan network
  - Before RAN, or
  - After DAE
- A reverse banyan network (RBN) can be used

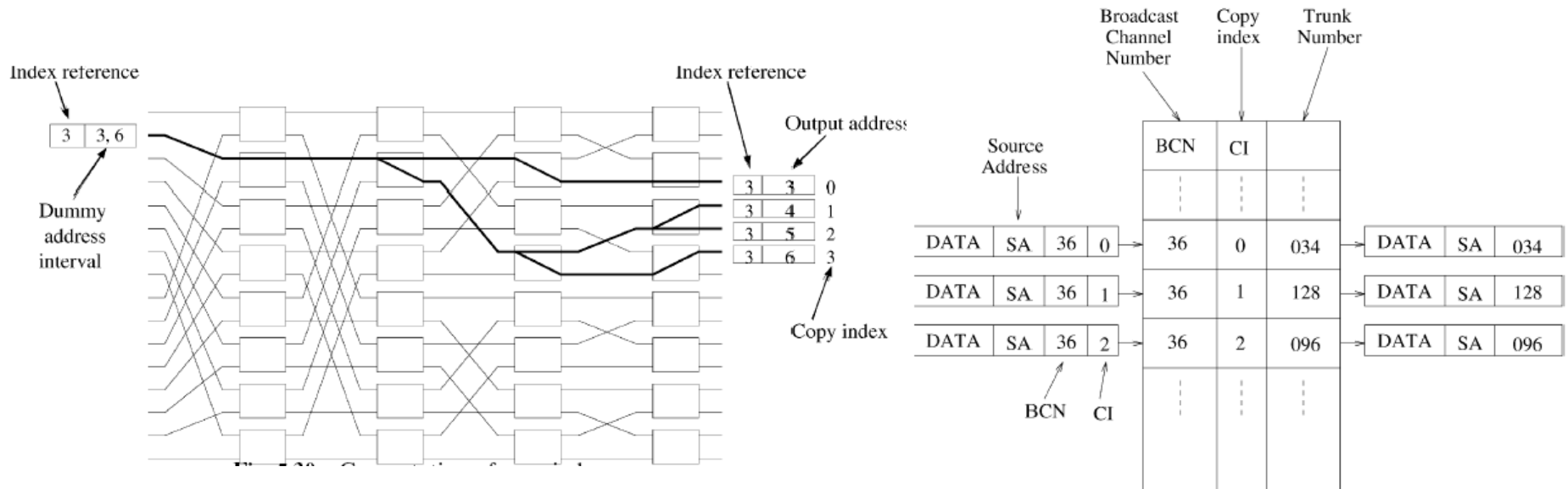




# MULTICAST COPY NETWORKS

## Decoding process

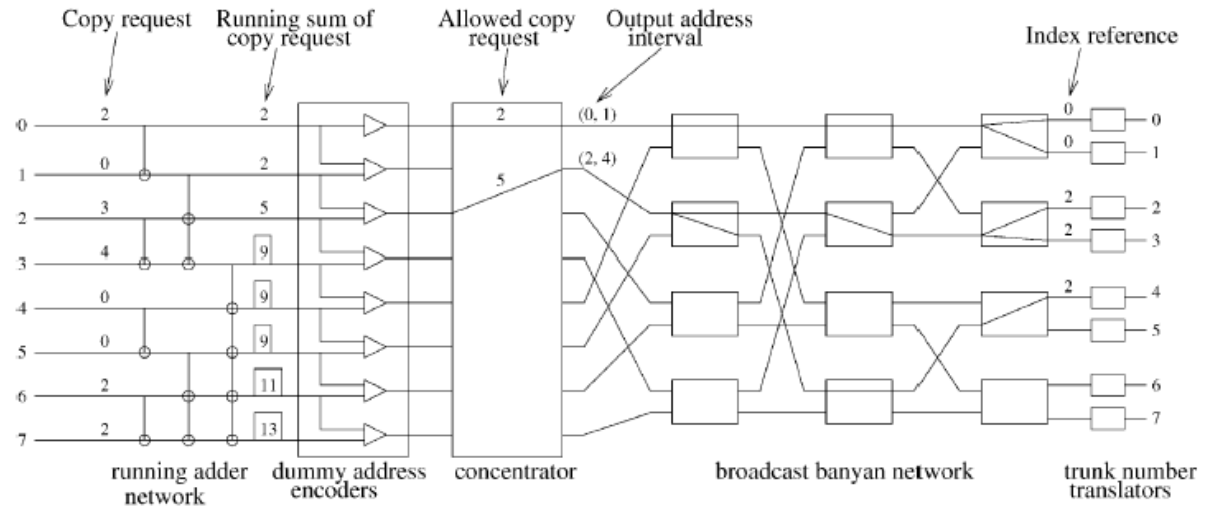
- When cell leaves broadcast banyan network,
  - The interval in the header is only one address
  - This address = min = max
  - Copy index = this address – index reference
  - TNT assigns the actual address to each copy
    - A simple table lookup
    - Search key: BCN and CI



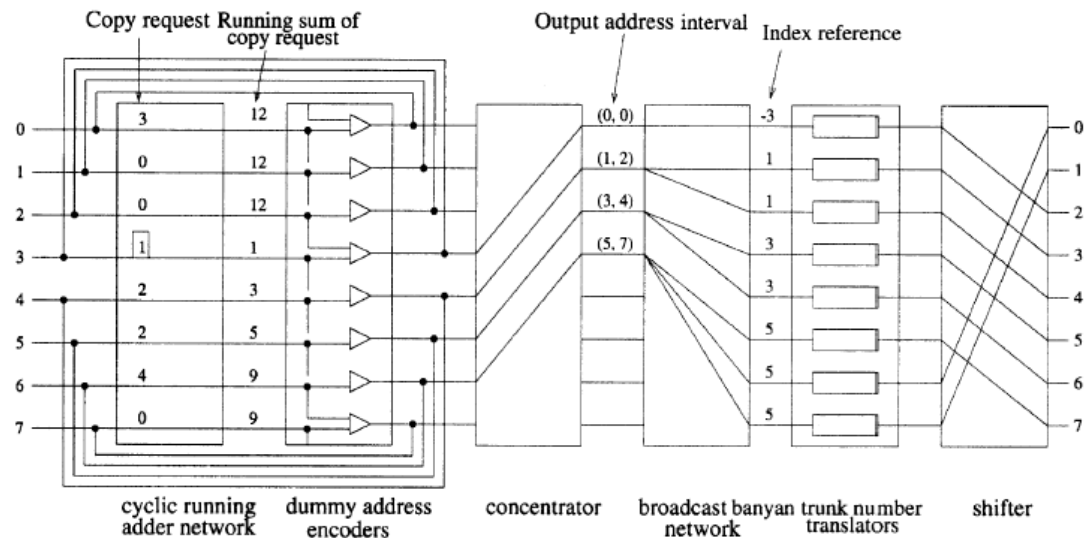
# MULTICAST COPY NETWORKS

## Overflow

- Copy network may not be able to do all copy requests
- An example
- Overflow problems
  - Performance
  - Unfairness



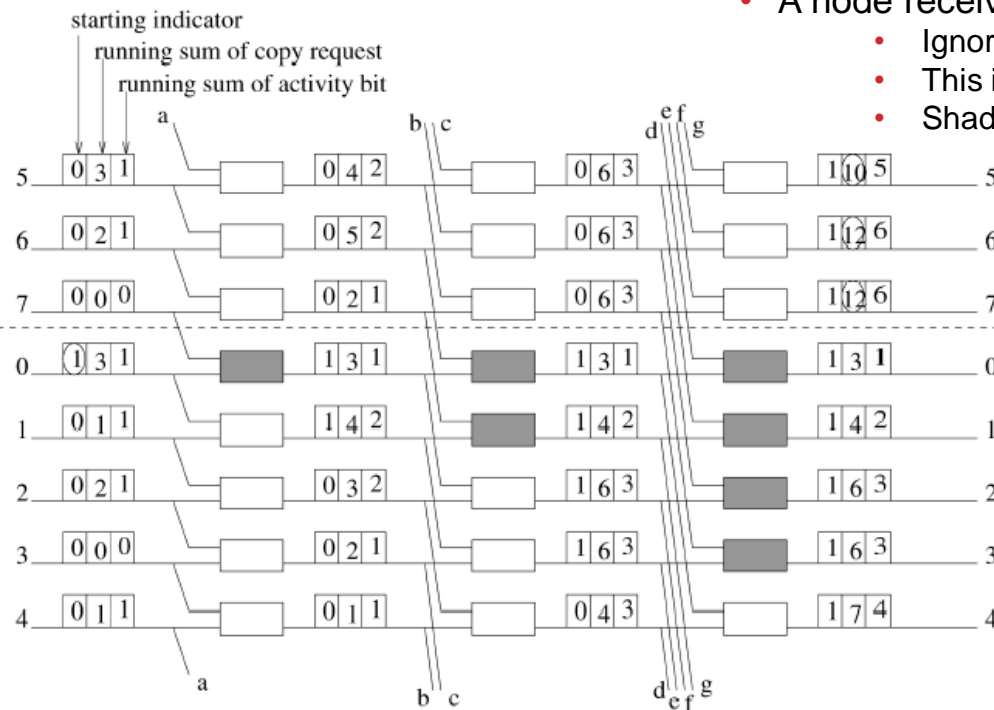
- Unfairness problem
  - Lower numbered inputs will have less overflow
  - Solution: CRAN instead of RAN
    - Adaptively changes RAN sum starting point



# MULTICAST COPY NETWORKS

## CRAN

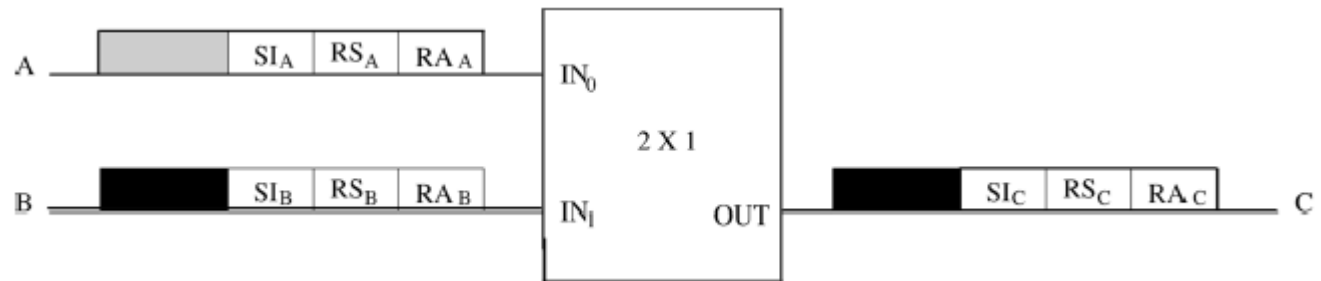
- Cell header fields
  - Starting indicator (SI)
  - Running sum (RS)
  - Routing address (RA)
- Initial values
  - SI: nonzero only for the starting point
  - RS: the number of copies
  - RA: 1 if port is active, otherwise 0
- At output, as the result:
  - RA has the running sum over activity bits
- A node receiving SI=1
  - Ignores its links
  - This is propagated
  - Shaded node in the figure



# MULTICAST COPY NETWORKS

## CRAN

- Header modification in a node



- The next starting point
  - No overflow  $\rightarrow$  same as the previous point
  - Overflow  $\rightarrow$  the first port facing the overflow
- RS updating

$$SI_0 = \begin{cases} 1 & \text{if } RS_{N-1} \leq N, \\ 0 & \text{otherwise.} \end{cases}$$

$$SI_i = \begin{cases} 1 & \text{if } RS_{i-1} \leq N \text{ and } RS_i > N, \\ 0 & \text{otherwise,,} \end{cases}$$

- SCN: starting copy number
  - Is sent back via feedback paths to input ports
  - So they know how many copies to serve

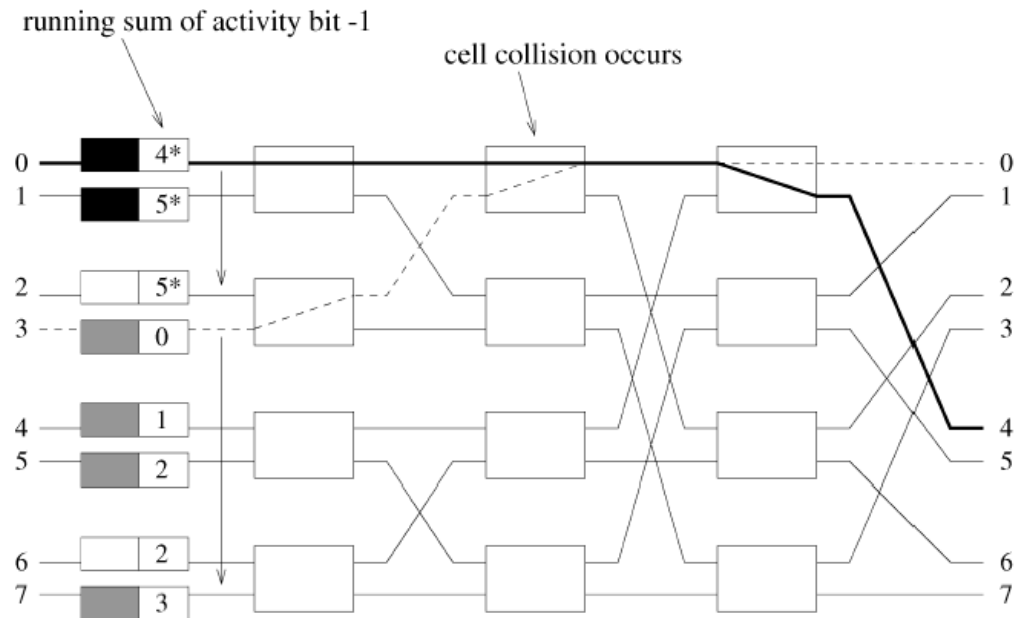
$$SCN_0 = RS_0,$$

$$SCN_i = \begin{cases} \min(N - RS_{i-1}, RS_i - RS_{i-1}) & \text{if } RS_{i-1} < N \\ 0 & \text{otherwise} \end{cases}$$

# MULTICAST COPY NETWORKS

## Concentration problem

- The starting point in CRAN may not be port 0
- Internal collisions may occur in reverse banyan network (RBN)



- Solution: an additional RAN before RBN

