

IN (Interconnection Network)

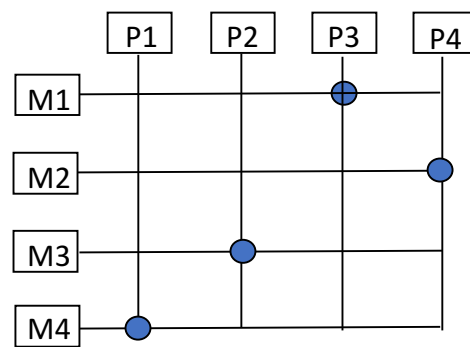
IN in MIMD shared-memory systems

Bus – shared among connected components;

flexible, low cost but, contention problem

→ not efficient for connecting many components (not scalable)

Cross-bar (x-bar) – switched network



ex) 4 simultaneous comm.

P1 – M4

P2 – M3

P3 – M1

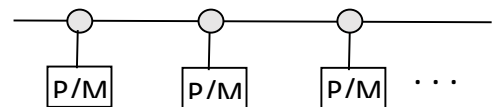
P4 – M2

Both bus and X-bar are dynamic topologies.

IN in MIMD distributed-memory systems

- Direct (static) IN – each switch is directly connected to a P-M pair;

ex) ring, mesh, torus, fully connected,

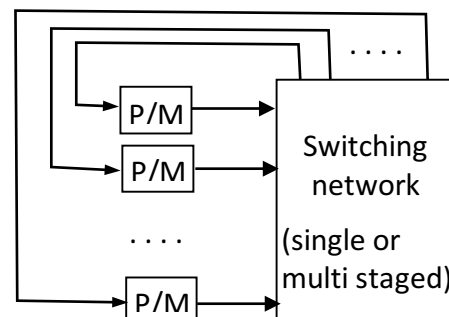


- Indirect (dynamic) IN – switches may not be directly connected to a node;

ex) X-bar (single stage),

Omega network (multi stages),

Beneš network (multi stages)



Static (direct) topologies

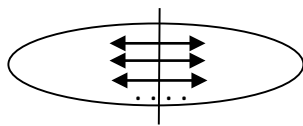
Assumptions/symbols used below:

network size (N) – # of nodes in IN

node degree (d) – # of links per internal node

network diameter (D) – any node can go to any other node within D hops
(i.e., at worst case (longest distance), best solution)

bisection width (connectivity) – # of simultaneous communications
between 2 halves of the network



(a) Completely connected

node degree $d = N-1$ ($N-1$ links/node)

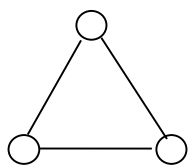
network diameter $D = 1$

fault tolerant

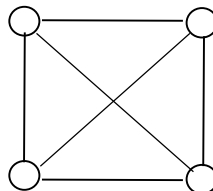
routing: directly connect to dest (1 hop) – high HW cost



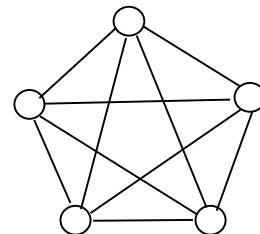
N=2



N=3



N=4



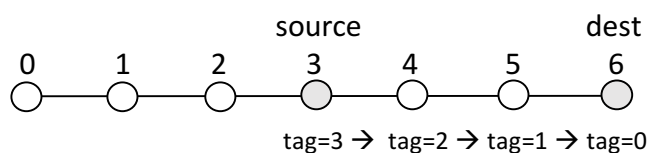
N=5

(b) Linear array

$d = 2$ links/node

$D = N-1$

not fault tolerant: one fault \rightarrow disconnects the entire list



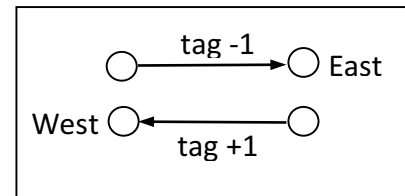
routing: compute 1-D tag value and perform E-W routing;

source tag = dest_id – source_id;

ex) source_id = 3, dest_id = 6

→ source tag = 6-3 = 3

trace until tag = 0;



(c) Ring (bidirectional)

d = 2 links/node

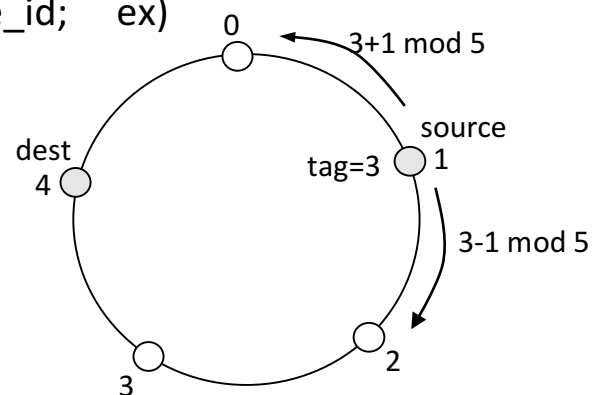
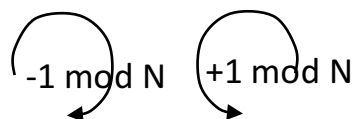
$D = \lfloor N/2 \rfloor$

fault tolerant (partly) – both directions available

Routing: source_tag = dest_id – source_id; ex)

new_tag = tag (+/- 1) MOD N;

trace until tag = 0;



(d) Mesh

2-D mesh

ex) $3^2=9$, $n=3$, $k=2$

$N = n^k$, where k, n are dimension, size of 1-D

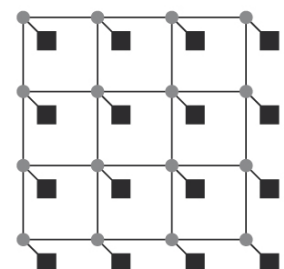
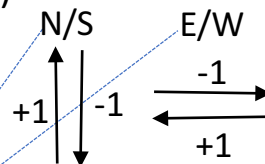
d = 2k = 4 links/node (inner nodes)

$D = k(n-1)$

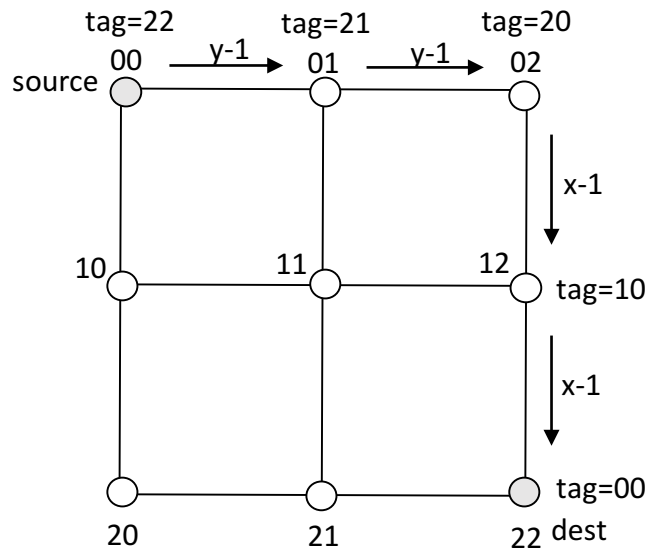
fault tolerant

Routing: X-Y routing, tag = X Y

source tag (2-D) = dest_id – source_id;



ex) source_id = 00, dest_id = 22; \rightarrow source tag = 22
 trace until tag = 00; if a node is busy, take any available path;



3-D mesh

ex) $3^3=27$, $n=3$, $k=3$

$N = n^k$, where k , n are dimension, size of 1-D

$d = 2k = 6$ links/node (inner nodes)

$D = k(n-1)$

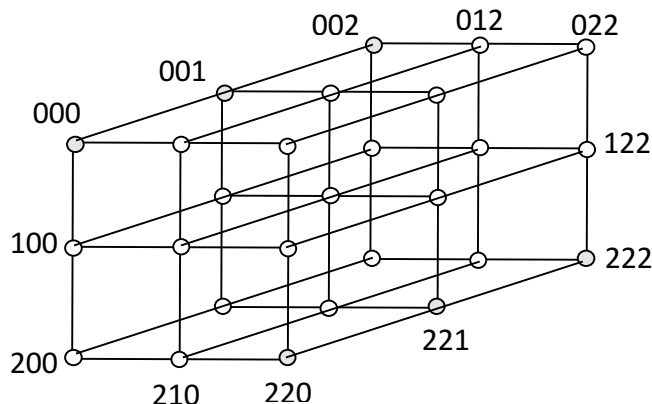
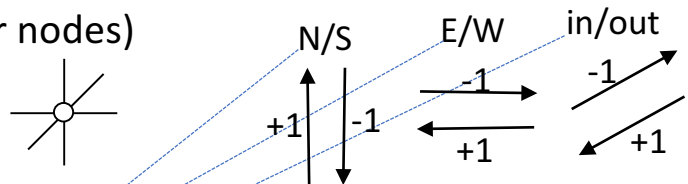
fault tolerant, expensive

Routing: X-Y-Z routing, tag = X Y Z

source tag value (3-D) = dest_id - source_id;

ex) source_id = 000, dest_id = 222; \rightarrow source tag = 222

trace until tag = 000; if a node is busy, take any available path;



ex) $N = 27$

$d = 2k = 2 \cdot 3 = 6$

$D = k(n-1) = 3(3-1) = 6$

(e) Torus (2-D and 3-D)

mesh + ring \rightarrow symmetric

$d = 2k$, where k is dimension, for all nodes

symmetric node degree, i.e.,

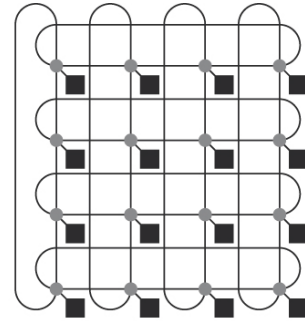
2-D torus: $d = 2k = 2 \times 2 = 4$ links/node

3-D torus: $d = 2k = 2 \times 3 = 6$ links/node

$D = k \lfloor \frac{\sqrt[k]{N}}{2} \rfloor$ //ex) 2-D 3x3 torus, $D = 2 \lfloor \frac{\sqrt[2]{9}}{2} \rfloor = 2 \times \lfloor 3/2 \rfloor = 2$

reduced network diameter, e.g., $N=9$, 2-D mesh ($D=4$) vs. 2-D torus ($D=2$)

3-D torus used in Cray T3D/TSE




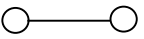
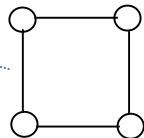
(f) Hypercube

network size $N = 2^k$, where k is dimension

node degree $d = k$

diameter $D = k$

dimension(k) # of nodes(N) node degree(d) diameter(D)

0 (0 cube)	1	0 link/node	0	
1 (1 cube)	2	1 link/node	1	
2 (2 cube)	4	2 links/node	2	
3 (3 cube)	8	3 links/node	3	
4 (4 cube)	16	4 links/node	4	

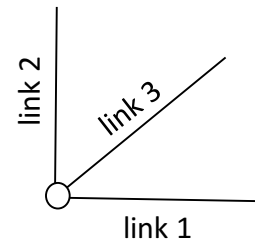
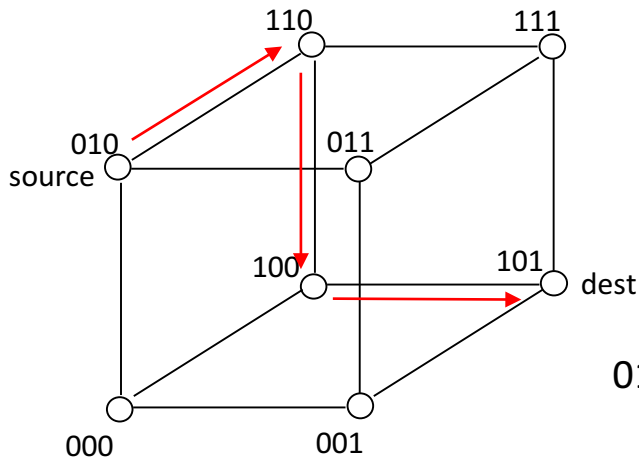
Routing: gray coding for node id – adjacent node_id differs in only 1 bit;

available links = source_id XOR dest_id;

repeat this with new source_id until no links are available;

possible policies: take left-most (or, right-most) available link first;

ex) 3-D cube with source_id = 010, dest_id = 101



$010 \text{ XOR } 101 \rightarrow 111$
 $\underbrace{\quad \quad \quad}_{\text{link3 link2 link1}}$
available links to take

if we use the policy of taking the left-most avail. one first,
 take link3: new source_id = 110, and perform
 $110 \text{ XOR } 101 \rightarrow 011$ (link2 and link1 are available);
 take link2: new source_id = 100, and perform
 $100 \text{ XOR } 101 \rightarrow 001$ (link1 is available);
 take link1: new source_id = 101, and perform
 $101 \text{ XOR } 101 \rightarrow 000$ (no links are available, stop)

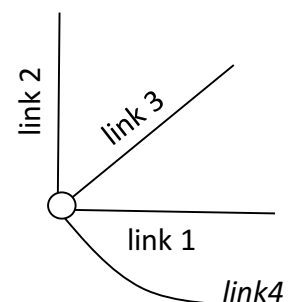
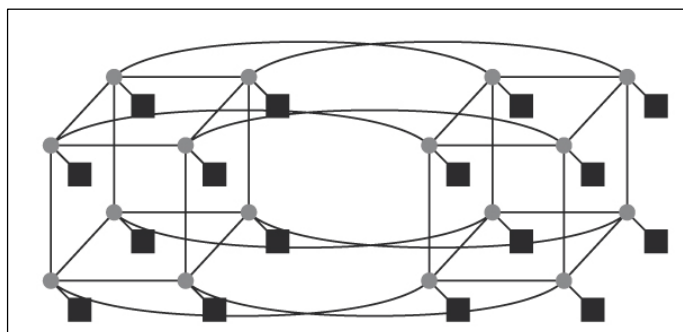
4-D cube routing:

gray coding and do analogous operation;

ex) source_id = 0011, dest_id = 1100

$0011 \text{ XOR } 1100 \rightarrow 1111$ (order of avail. link4,3,2,1)

4-D hypercube

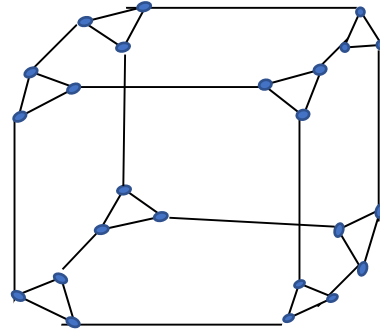


(g) Cube-connected cycles (CCC)

3-D CCC (3-CCC)

k-CCC : $k \cdot 2^k$ nodes

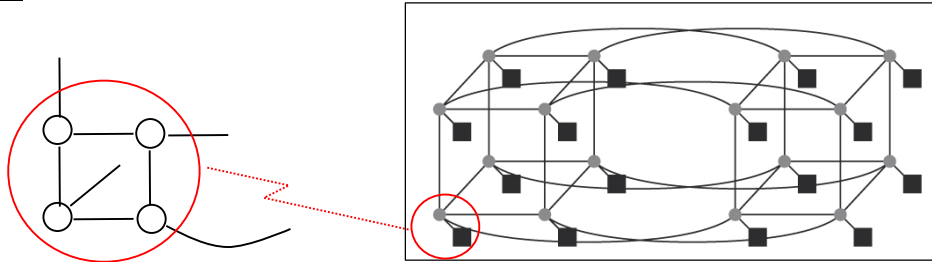
ex) 3-CCC,
 $3 \cdot 2^3 = 24$ nodes



fixed node degree (independent from k), $d=3$

network diameter $D = 2k - 1 + \lfloor k/2 \rfloor \approx 2k$

4-CCC : refined 4 nodes in each node of 4-cube



$d = 3$ (fixed)

$D = 2k - 1 + \lfloor k/2 \rfloor = 2 \cdot 4 - 1 + \lfloor 4/2 \rfloor = 9$

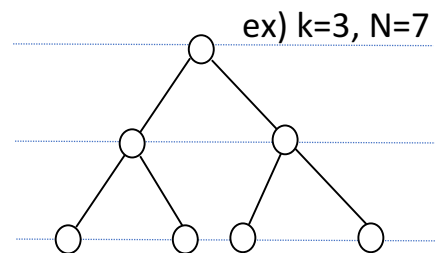
(h) binary tree

node degree = 3 (3 links/node), diameter $D = 2(k - 1)$

k-level binary tree: $N = \frac{2^k - 1}{2 - 1}$ nodes

k-level m-way tree: $N = \frac{m^k - 1}{m - 1}$ nodes

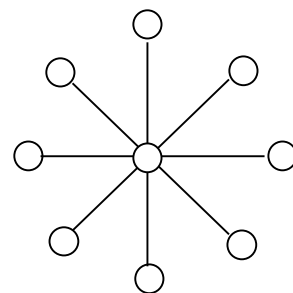
ex) $k=3 \rightarrow N=7$ nodes (Fig(h))



(k) star

node degree (internal node) $d = N - 1$

diameter $D = 2$



Dynamic IN

Bus – not used in SIMD, used in shared-mem MIMD

X-bar (crossbar) – single stage – shared-mem (P-M, P-P), dist-mem(P/M-P/M)

MIN (multi-stage IN) $\left[\begin{array}{l} \text{Omega network} \\ \text{Bene\hat{s} network} \end{array} \right)$ based on multi-stage shuffling

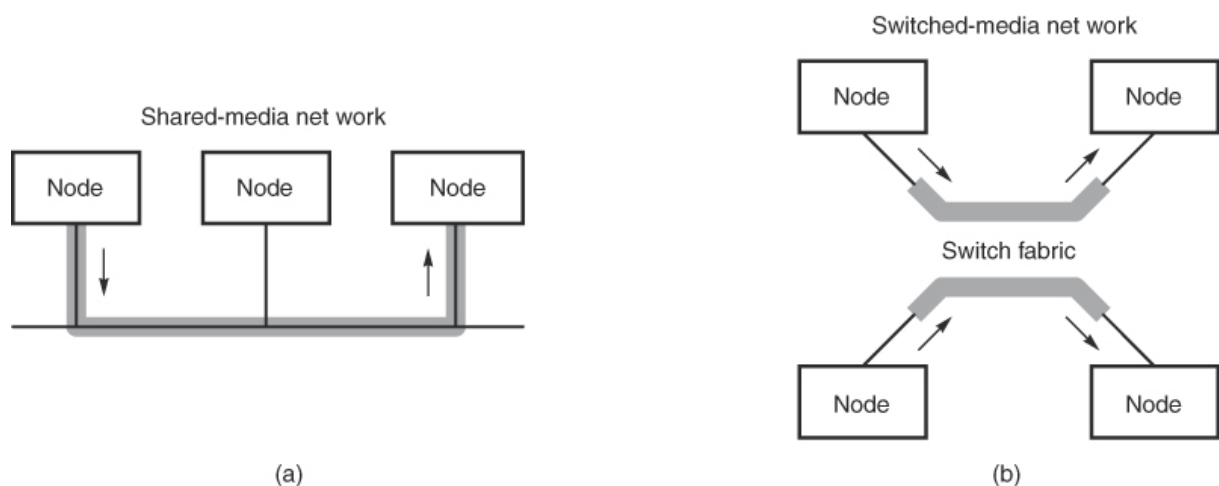


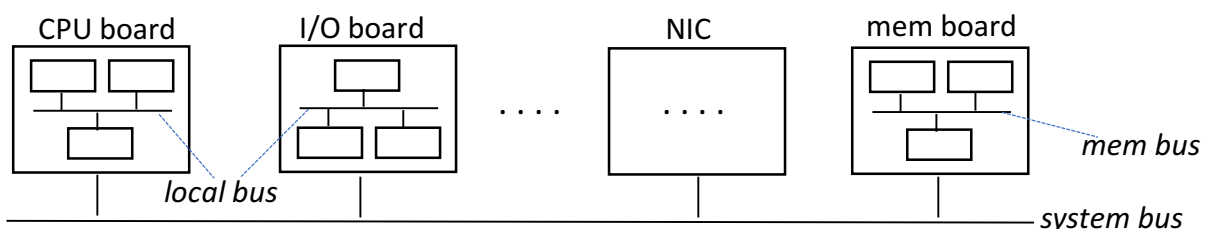
Figure F.8 (a) A shared-media network versus (b) a switched-media network. Ethernet was originally a shared media network, but switched Ethernet is now available. All nodes on the shared-media must dynamically share the raw bandwidth of one link, but switched-media networks can support multiple links, providing higher raw aggregate bandwidth.

Bus: a collection of wires/connectors for data transactions among processors, memory modules and peripheral devices.

system bus (P \leftrightarrow M) – datapath, address/control lines;

local bus – in each board (CPU board, mem board, I/O board);

memory bus – local bus in the memory board;



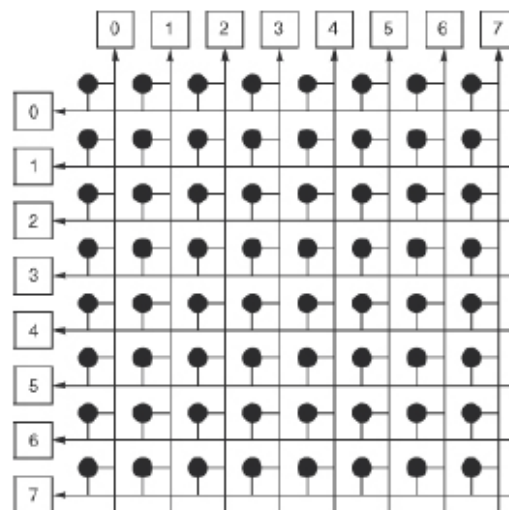
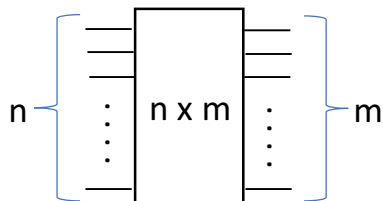
disadvantages: time shared use: bus access – one at a time → contention;
 limited scalability;
 processor bandwidth is a portion of the total bus bandwidth;

Crossbar switches

much higher bandwidth than bus;
 single stage switched network via cross point on/off;
 number of cross points → HW cost/complexity;
 usage: P-P communication in SMP, MPP, COW;
 P-M communication in SMP, PVP;

ex) 8x8 X-bar switch

in general, $n \times m$, where
 $n, m = 2^x$

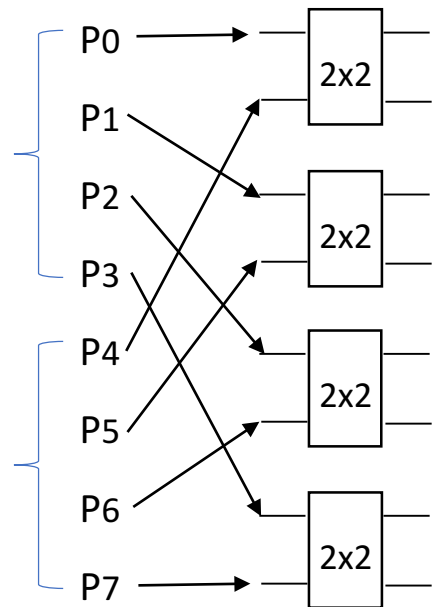


The crossbar network requires N^2 crosspoint switches, shown as black dots

Concurrent communication is available, e.g., PE1 – M2 and PE3 – M1;
 uniform latency – only 1 set of switches in any path;
 contention case → one waits;
 operation of MIMD requires P-M IN be changed in a dynamic fashion;

Shuffling

perfect shuffle

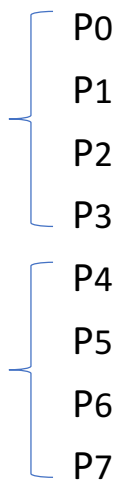


Concept:

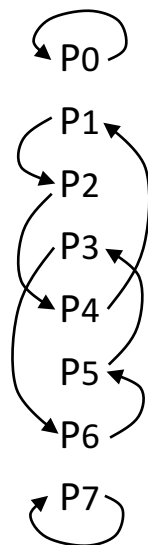
divide PEs into 2 groups, e.g.,
(0,1,2,3) and (4,5,6,7);

2 groups are then merged, s.t.,
0 is adjacent to 4;
1 is adjacent to 5;
2 is adjacent to 6;
3 is adjacent to 7;

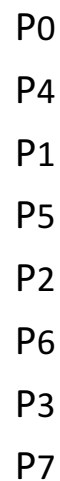
before shuffle



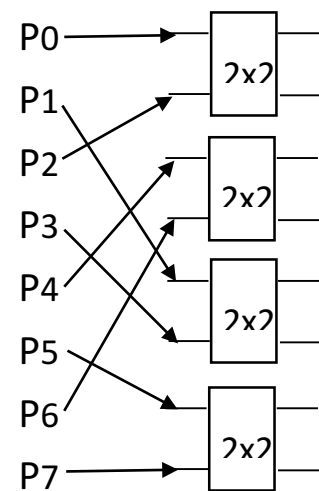
shuffle



after shuffle



reverse perfect shuffle



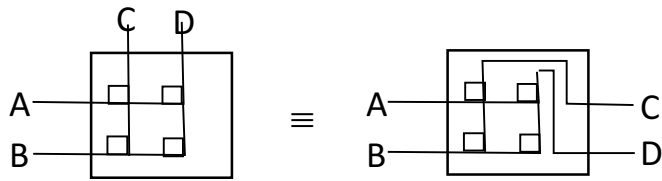
if $i < N/2$, shuffle $(i) = 2*i$

if $i \geq N/2$, shuffle $(i) = [(2*i) \text{ MOD } N] + 1$

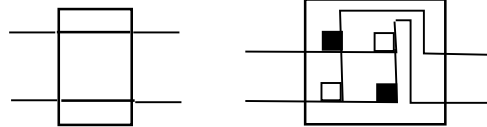
node_id: $\overbrace{x \ y \ z}^{\text{shuffle}}$

ex) P0 (000) → 000; P1 (001) → 010; P2 (010) → 100

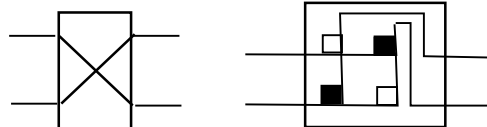
2x2 X-bar switch operations



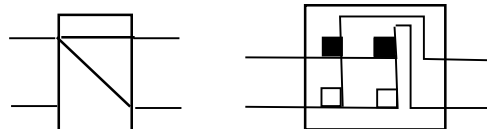
pass:



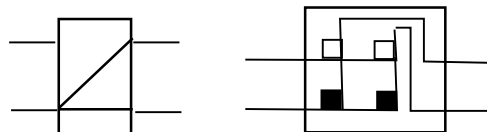
cross:



upper-
broadcast:



lower-
broadcast:



MIN (multi-stage IN)

dynamic switch connection (on/off);

fixed inter-stage connections between adjacent stages;

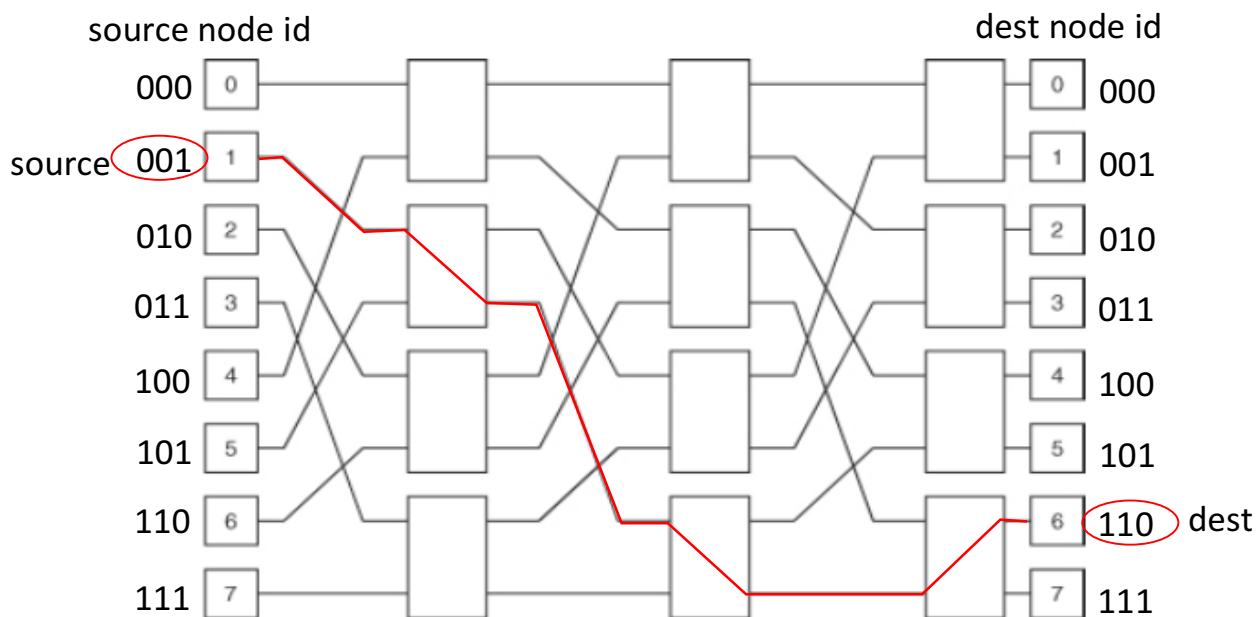
ex) Omega network, Beneš network

Omega network

N x N Omega network with 2x2 X-bar switches:

- $\log N$ stages, each stage with $\frac{N}{2}$ 2x2 X-bar switches
→ total $(\log N) * \frac{N}{2}$ 2x2 X-bar switches;
- perfect shuffle between stages;
- packet switching (data, dest-address);
read in each switch

ex) N=8 (i.e., 8x8)



routing: dest_id routing

ex) dest_id = 1 1 0

1st switch low 2nd switch low 3rd switch high

- Any node to any node comm. is possible, but \exists only 1 unique path from an input (source) to an output (dest).

problem: blocked case; → Sol: buffered in switch and wait

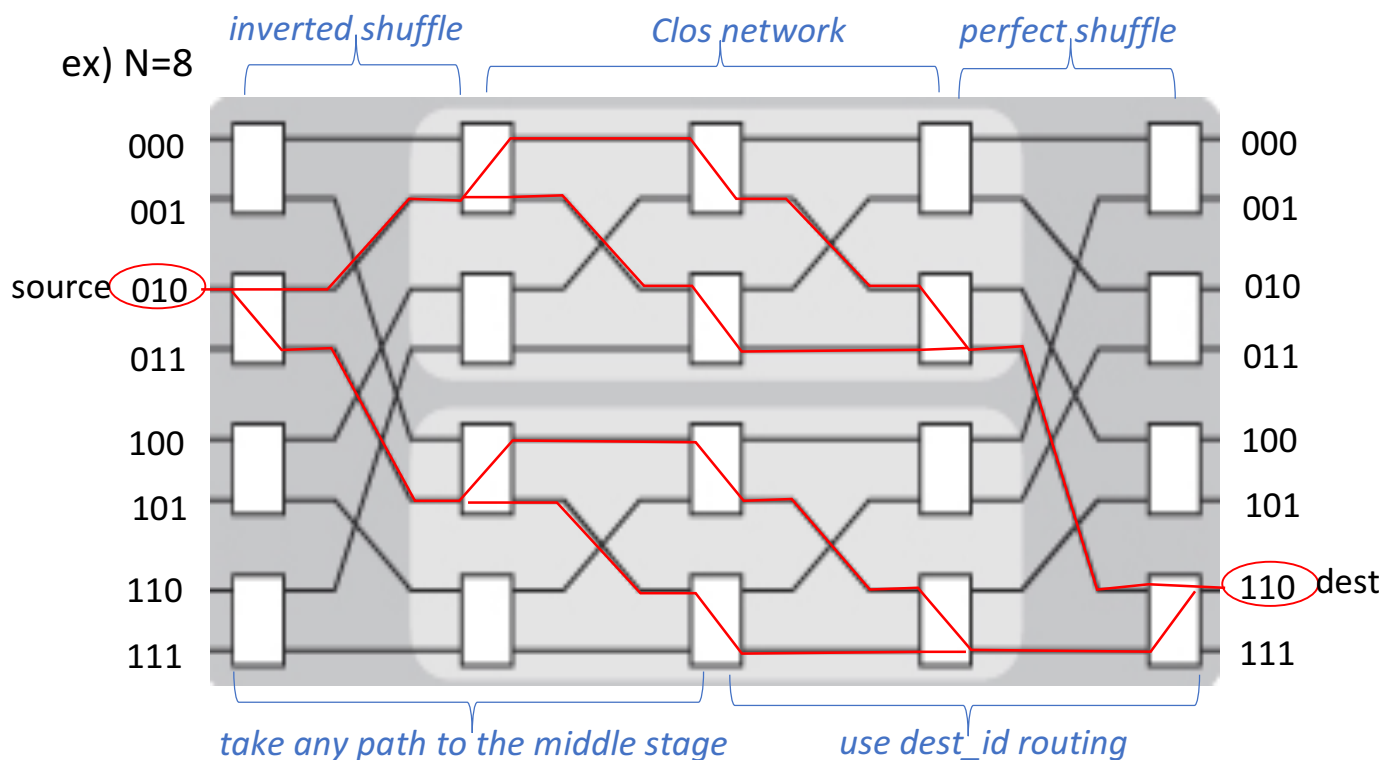
Beneš network

- \exists multiple paths between a source and a dest.;
 \rightarrow reduced blocking, increased bandwidth;
- More complex HW;
- increased network latency;
 increased # of stages \rightarrow complexity in routing path computation;
 so, suitable for circuit switching since the routing path computation takes a considerable amount of time.

of stages = $2(\log N) - 1$

ex) $N=8 \rightarrow 5$ stages, $N=16 \rightarrow 7$ stages

total # of 2x2 switches = $\frac{N}{2} [2(\log N) - 1]$



total # of paths between a source and dest = $2^{\lfloor \text{total \# of stages} / 2 \rfloor}$

ex) $N=8 \rightarrow 2^2 = 4$ paths

N=16 (16x16) Beneš network

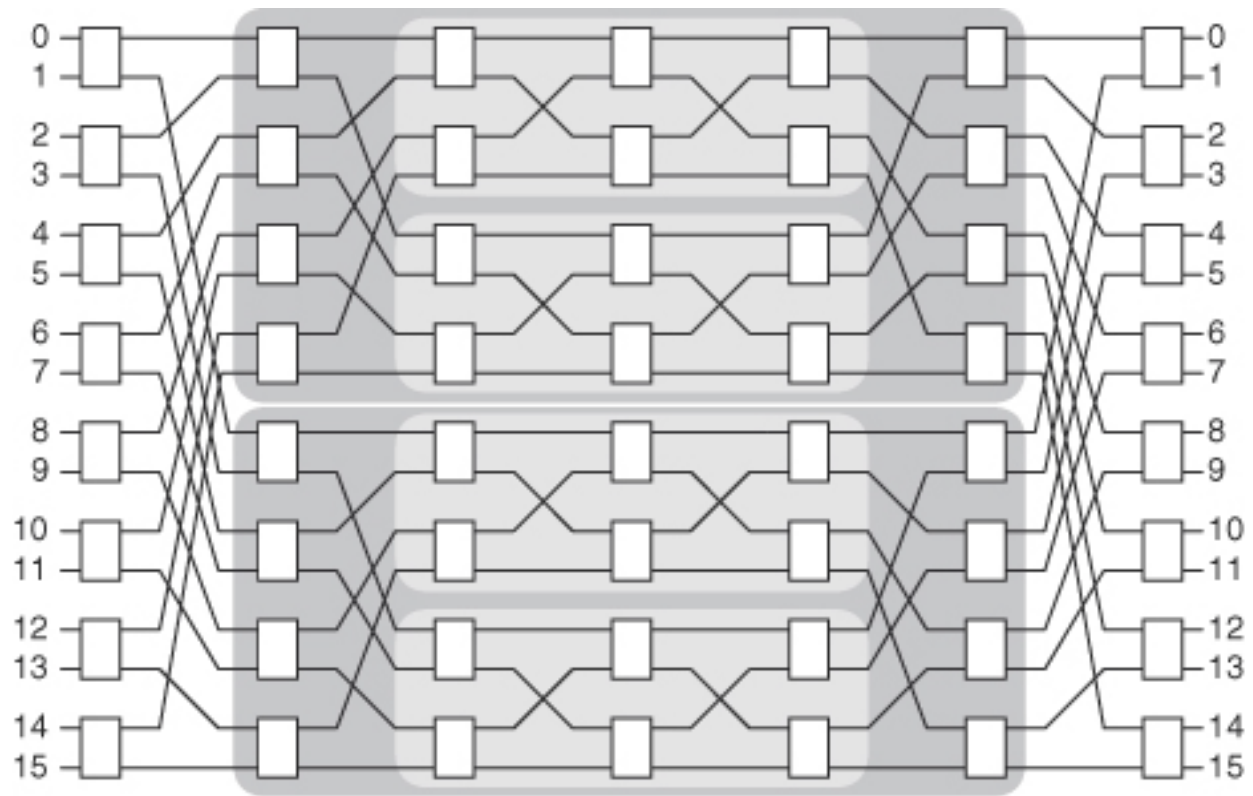


Figure F.12 Beneš network. A 16-port Clos topology, where the middle-stage switches shown in the darker shading are implemented with another Clos network whose middle-stage switches shown in the lighter shading are implemented with yet another Clos network, and so on, until a Beneš network is produced that uses only 2 x 2 switches everywhere.

$$2(\log N) - 1 = 7 \text{ stages}$$

$$\text{total \# of paths between a source and a dest} = 2^{\lfloor \text{total \# of stages} / 2 \rfloor} = 2^3 = 8$$

$$\text{total \# of 2x2 switches} = \frac{N}{2} * [(2 \log N) - 1] = 8 * 7 = 56$$

Mk book *distributed-memory indirect IN (dynamic)*

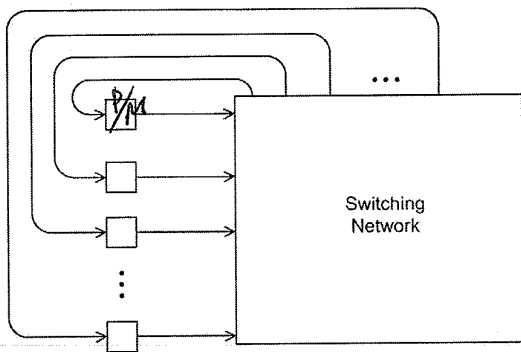


Fig 2.13 A generic indirect network

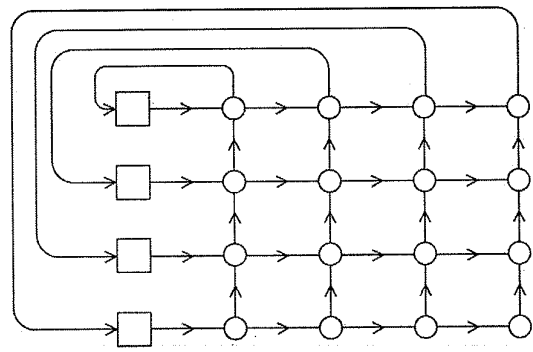


Fig 2.14 A crossbar interconnect for distributed-memory
(4x4) dynamic single stage

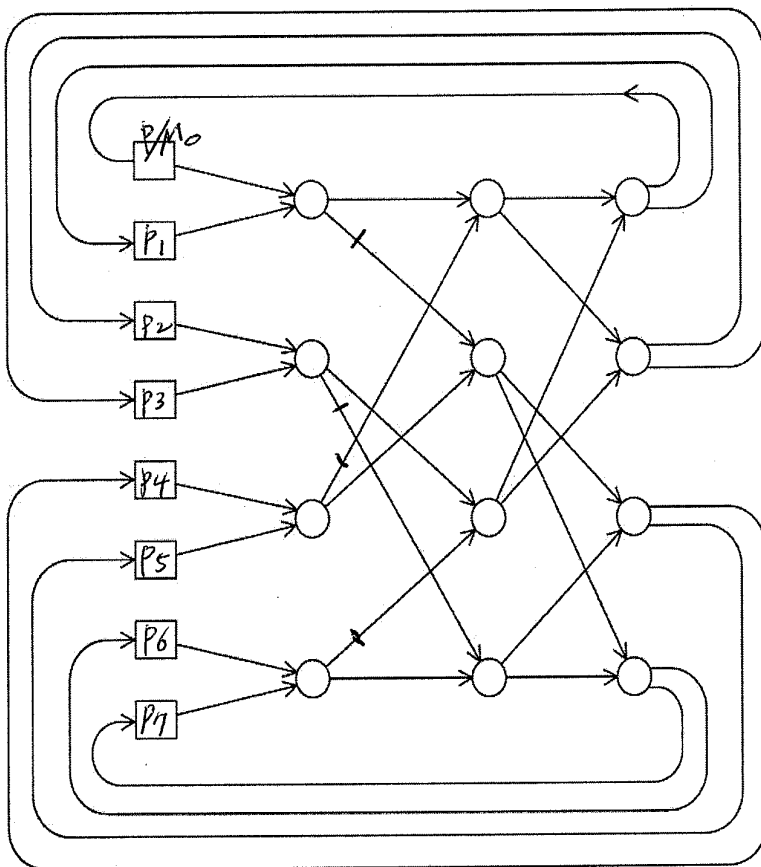


Fig 2.15 An omega network

dynamic multi-stage

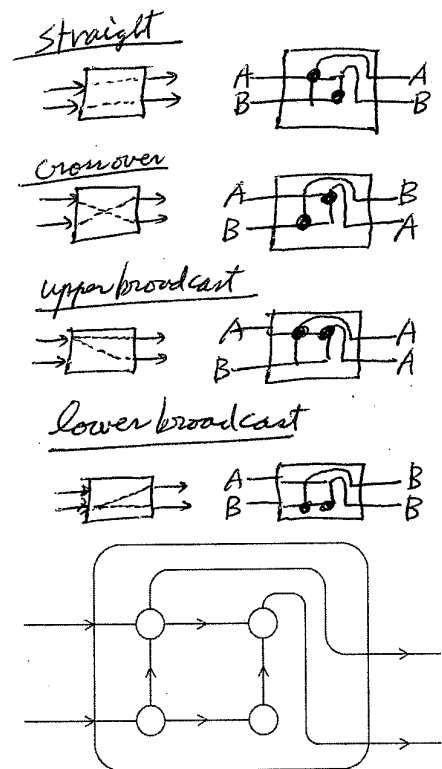


Fig 2.16 A switch in an omega network
(2x2)