

## 第 2 章 相互結合網

### 2 . 1 相互結合網の分類

#### 2 . 1 . 1 分類項目

##### ( 1 ) 制御方式

集中制御方式

分散制御方式

##### ( 2 ) 交換方式

パケット交換方式 ( packet switching ) :

手紙による通信

- ・ 蓄積交換 (store and forward)
- ・ ワームホール (worm hole)
- ・ バーチャルカットスルー  
(virtual cut through)

回線交換方式 (circuit switching) :

電話による通信

### ( 3 ) トポロジ

静的網 : 直接結合 (direct connection) 網

完全網を規則的に簡略化

動的網 : 間接結合 (Indirect Connection) 網

クロスバスイッチを規則的に簡略化

## 2.1.2 評価項目

- (1) 距離 (distance)
- (2) 次数 (degree)
- (3) 総スイッチ数 / 総リンク数
- (4) 拡張性 (scalability)
- (5) 3次元実装の容易性
- (6) 耐故障性 (fault tolerance)
- (7) 多様な網の埋込み能力 (embedability)

静的網の場合：

多数の他の網の埋込みが可能である必要  
並列処理の応用プログラム  
論理構造に適合した特定のトポロジ上で  
最も高速に実行

( 8 ) ルーティングの容易性



## 2.1.3 基本通信パターン

送信ノード番号  $X$  の 2 進表示:  $(a_n, \dots, a_2, a_1)$

受信ノード番号  $Y$  の 2 進表示:  $(b_n, \dots, b_2, b_1)$

通信パターンの数 (受信ノードにダブリなし):  $N!$

送受信ノードの対応関係: ひとつの置換

送信

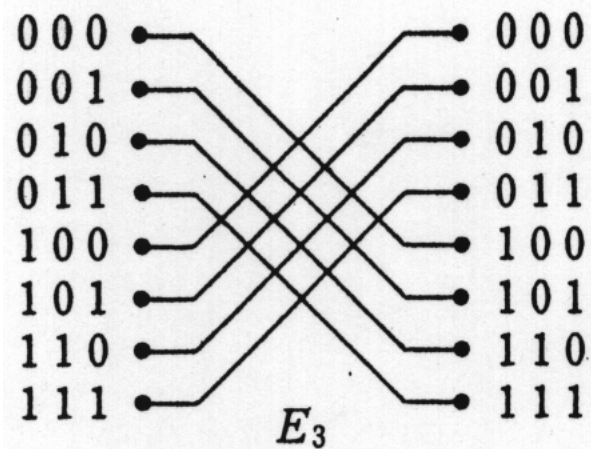
受信

$(1, 2, 3, 4) \leftrightarrow (2, 3, 4, 1)$

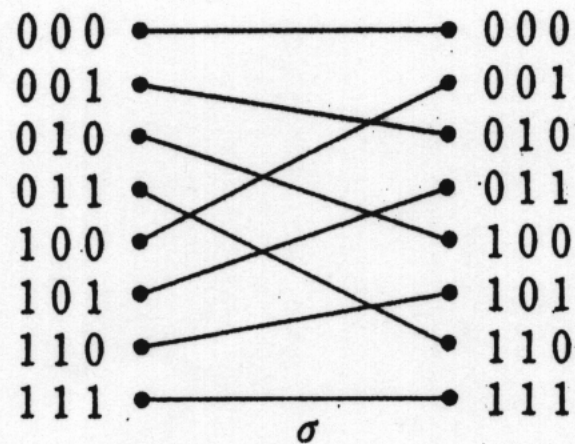
(1) エクスチェンジ置換

$$y = E_i(x) = (a_n, \dots, \sim a_i, \dots, a_1)$$

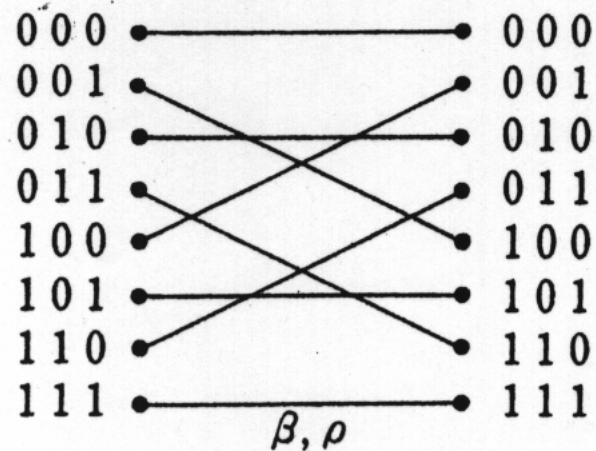
$\sim a_i$  は  $a_i$  の否定を表す。



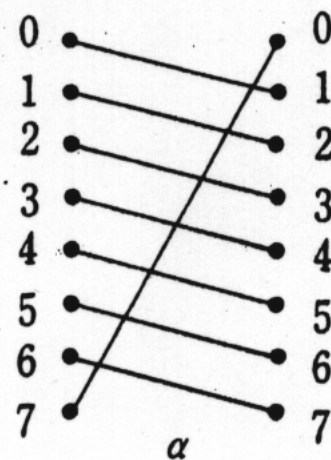
(a) エクスチェンジ置換



(b) シャフル置換



(c) バタフライ置換と  
ビット逆転置換



(d) シフト置換

## ( 2 ) シャフル置換

$$y = \sigma(x) = (a_{n-1}, a_{n-2}, \dots, a_1, a_n)$$

k - サブシャフル  $_k$  、

k - スーパシャフル  $^k$

$$y = \sigma_k(x) = (a_n, \dots, a_{k+1}, a_{k-1}, \dots, a_1, a_k)$$

$$y = \sigma^k(x) = (a_{n-1}, \dots, a_{n-k+1}, a_n, a_{n-k}, \dots, a_1)$$

## ( 3 ) バタフライ置換

$$y = \beta(x) = (a_1, a_{n-1}, \dots, a_2, a_n)$$

k - サブバタフライ  $_k$  、

k - スーパーバタフライ <sup>k</sup>

$$y = \beta_k(x) = (a_n, \dots, a_{k+1}, a_1, a_{k-1}, \dots, a_2, a_k)$$

$$y = \beta^k(x) = (a_{n-k+1}, a_{n-1}, \dots, a_{n-k+2}, a_n, a_{n-k}, \dots, a_1)$$

( 4 ) ビット逆転置換

$$y = \rho(x) = (a_1, a_2, \dots, a_n)$$

k - サブビット逆転置換 <sup>k、</sup>

k - スーパービット逆転置換 <sup>k</sup>

$$y = \rho_k(x) = (a_n, \dots, a_{k+1}, a_1, a_2, \dots, a_{k-1}, a_k)$$

$$y = \rho^k(x) = (a_{n-k+1}, a_{n-k+2}, \dots, a_{n-1}, a_n, a_{n-k}, \dots, a_1)$$

## ( 5 ) シフト置換

$$y = \alpha(x) = |x+1|2^n$$

k - サブシフト置換  $\alpha_k$  、

k - スーパシフト置換  $\alpha^k$

$$y = \alpha_k(x) = |x+1|2^k + \lfloor x/2^k \rfloor 2^k$$

$$y = \alpha^k(x) = |x+2^{n-k}|2^n$$

$||$  : モジュロ計算、

$\lfloor a \rfloor$  : 床関数 ( floor function )

恒等置換  $i$

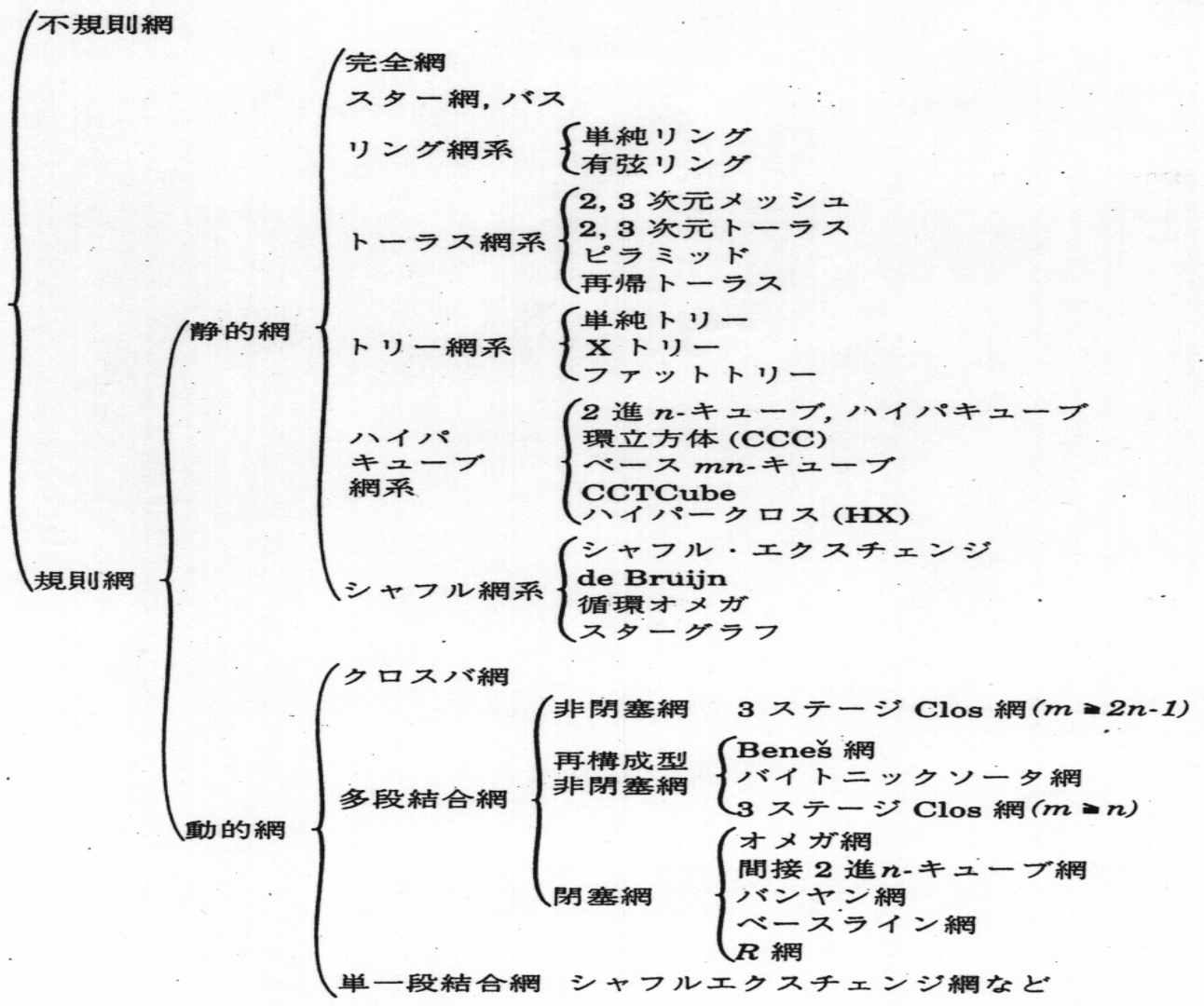
$$i(x) = x$$

逆置換  $\pi^{-1}$

$$\pi\pi^{-1} = \pi^{-1}\pi = i$$

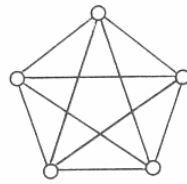
逆シャフル置換

$$y = \sigma^{-1}(x) = (a_1, a_n, a_{n-1}, a_{n-2}, \dots, a_2)$$

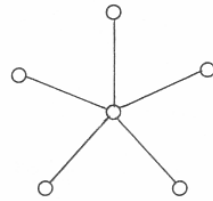


## 7.2 相互結合網

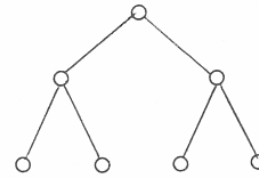
静的網



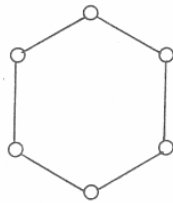
(a) 完全網



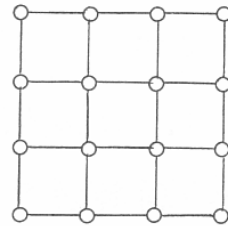
(b) スター網



(c) 木状網

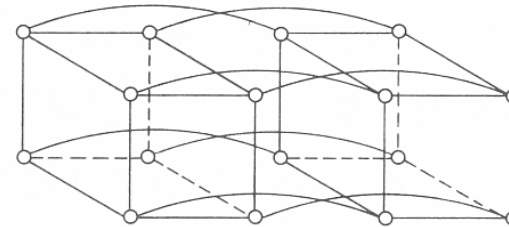


(d) リング網



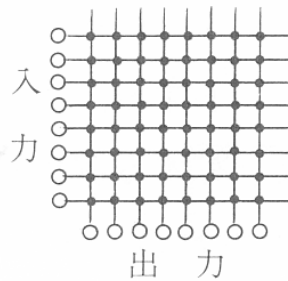
(e) 格子網(トーラス網)

(最上・下, 最左・右を結合)



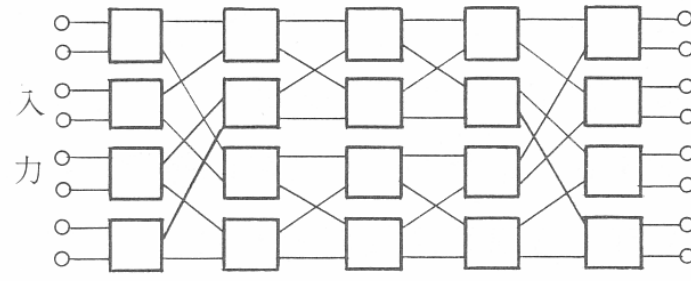
(f) ハイパーキューブ網

動的網



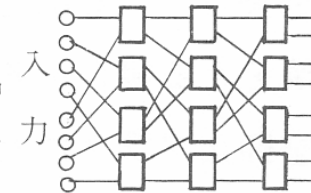
・: スイッチ

(a) クロスバー網



□: 2×2 スイッチ

(b) Beneš 網



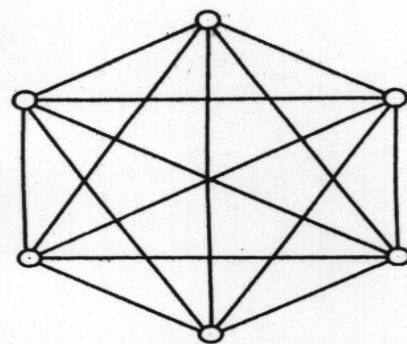
□: 2 入力 2 出力  
交換スイッチ

(c) オメガ網

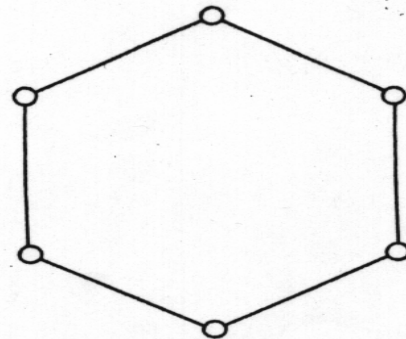
多段結合網

○: 演算装置やプロセッサ

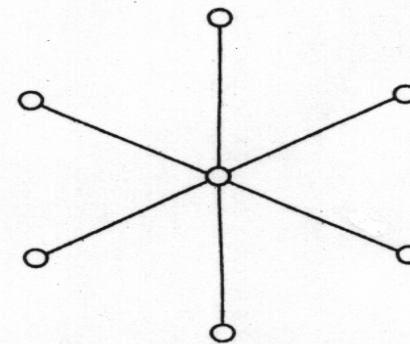




完全網

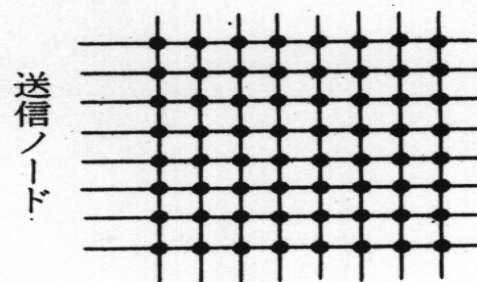


リング網

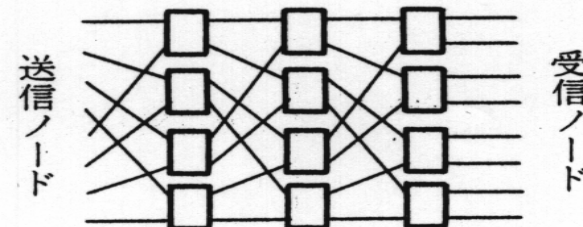


スター網

(a) 静的網



受信ノード  
● スイッチ  
クロスバ網



□ 2入力-2出力  
交換スイッチ

オメガ網

(b) 動的網

図 2.2 静的網と動的網

## 2 . 2 静的網

### 2 . 2 . 1 分類

( 1 ) 完全網

( 2 ) スター網

( 3 ) リング網系

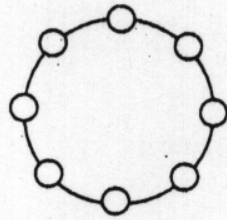
( 4 ) トーラス網系

2次元メッシュ網

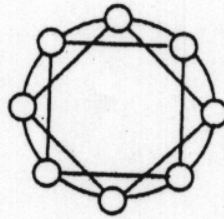
2次元トーラス網

3次元トーラス網

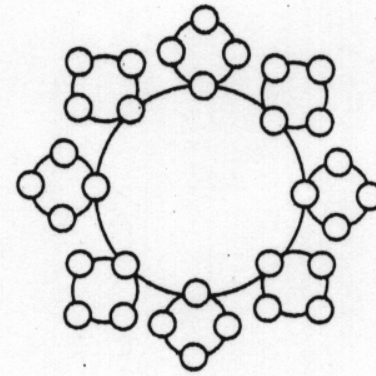
ピラミッド網



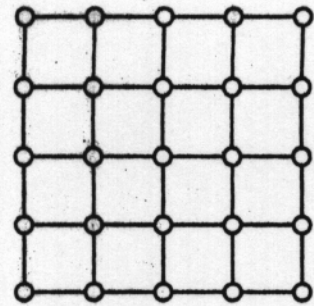
(a).単純リング



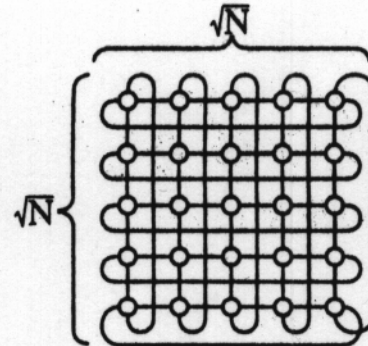
(b) 有弦リング



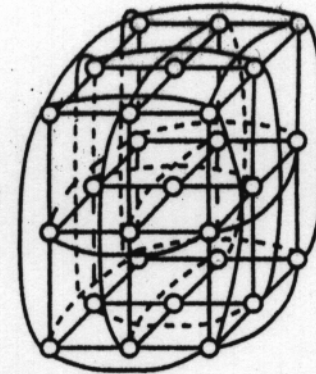
(c) 階層リング



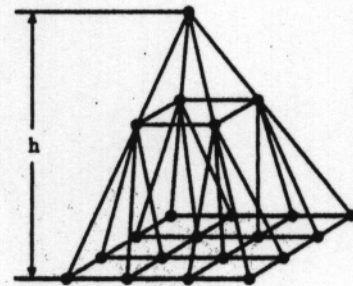
(a) 2次元メッシュ



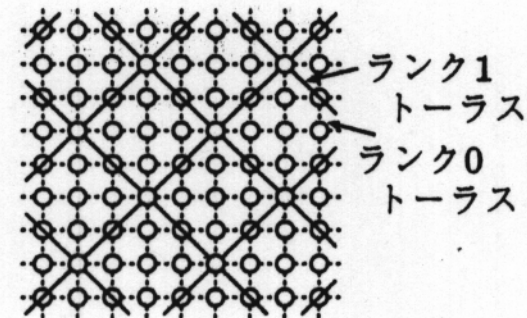
(b) 2次元トーラス



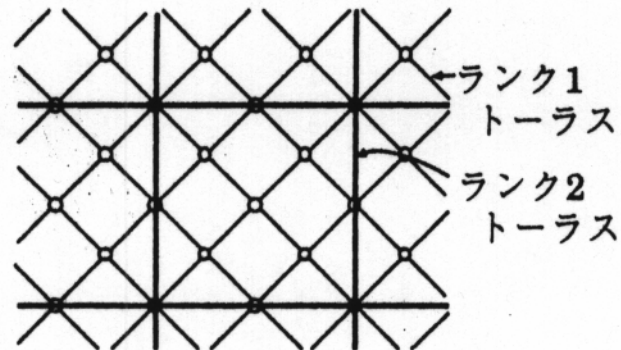
(c) 3次元トーラス



(d) ピラミッド



(i) ランク1トーラス



(ii) ランク2トーラス

(e) 再帰トーラス

再帰トーラス網

(RDT、Recursive Diagonal Torus)

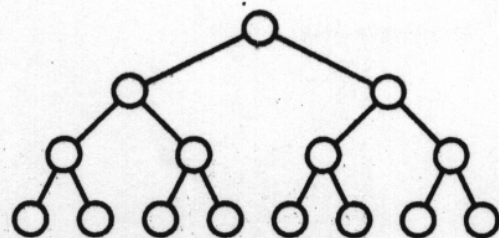
直径は $N=1024$ で 7、 $4096$ で 8、 $16384$ で 10 程度

( 5 ) トリー網系

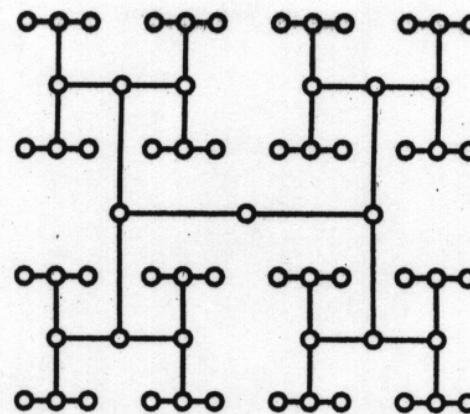
単純トリー網

X トリー

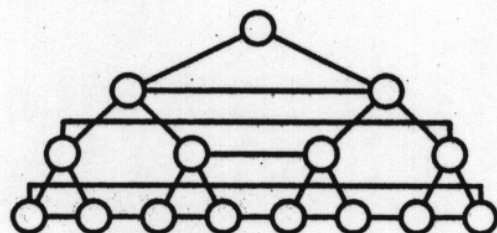
ファットトリー ( fat tree )



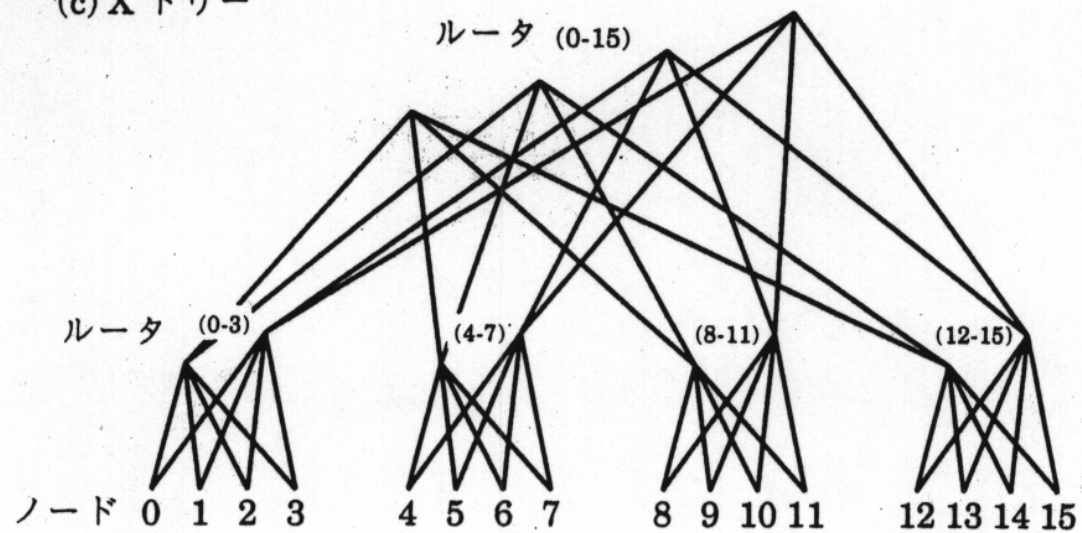
(a) 単純トリー



(b) トリーの VLSI チップへの埋込み



(c) X トリー



(d) ファットトリー (CM-5)

## ( 6 ) ハイパーキューブ網系

### ハイパーキューブ網

#### 特徴

ルーティング容易

拡張可能性

多様な網の埋め込み

リング：ノード番号を反射 2 進で

トーラス：ノード番号 2 次元反射 2 進で

トリー： $2^n$  の 1 つをダミーノード

# 結合

$(a_n, a_{n-1}, \dots, a_2, a_1)$



$(a_n, a_{n-1}, \dots, a_2, \bar{a}_1)$

$(a_n, a_{n-1}, \dots, \bar{a}_2, a_1)$

$\vdots$

$(a_n, \bar{a}_{n-1}, \dots, a_2, a_1)$

$(\bar{a}_n, a_{n-1}, \dots, a_2, a_1)$

ハミング距離: 1

$(001) \quad (110)$

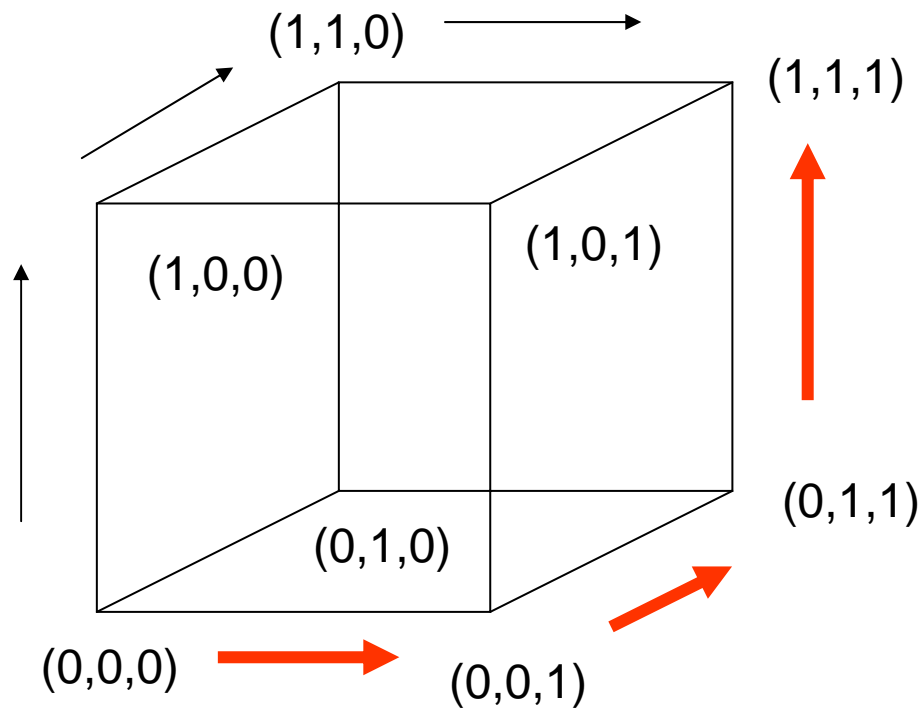
ハミング距離: 3

ハミング距離: 各桁の比較。異なる桁の数

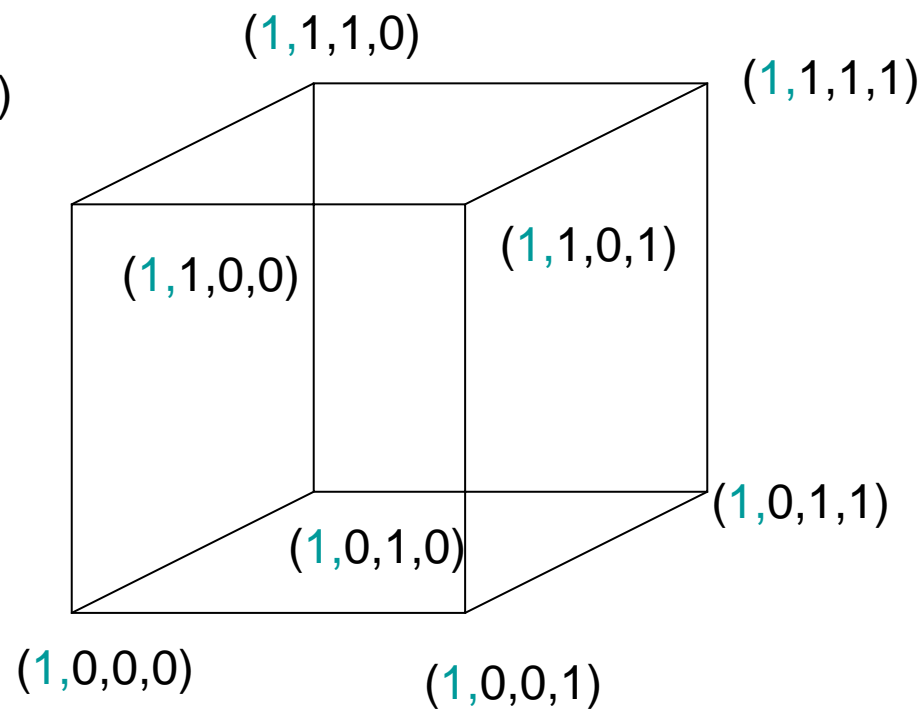
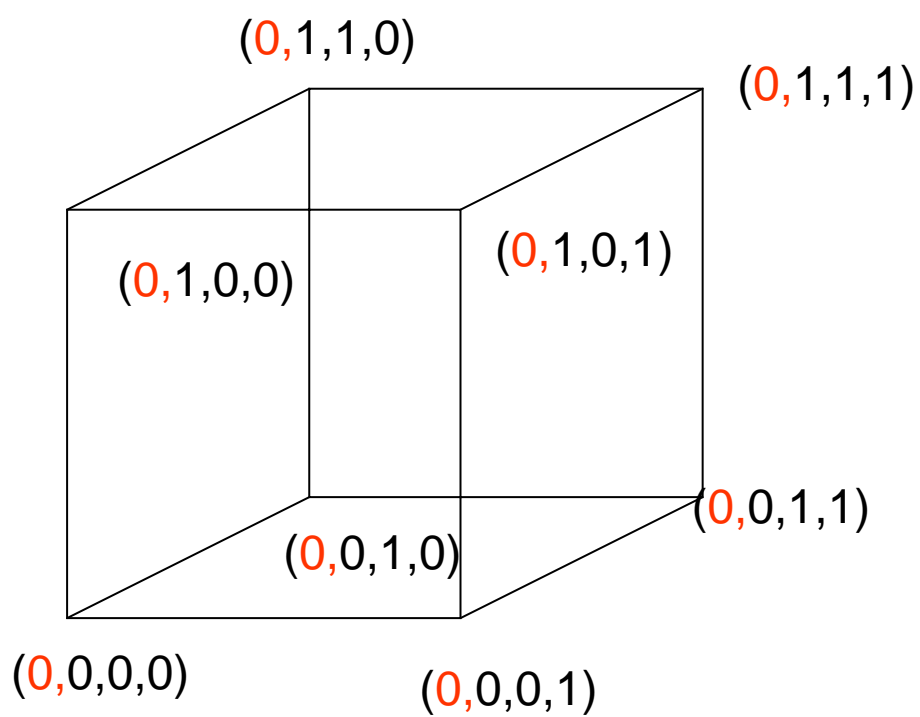


# ルーティング

## 下位ビットから宛先に合わせる

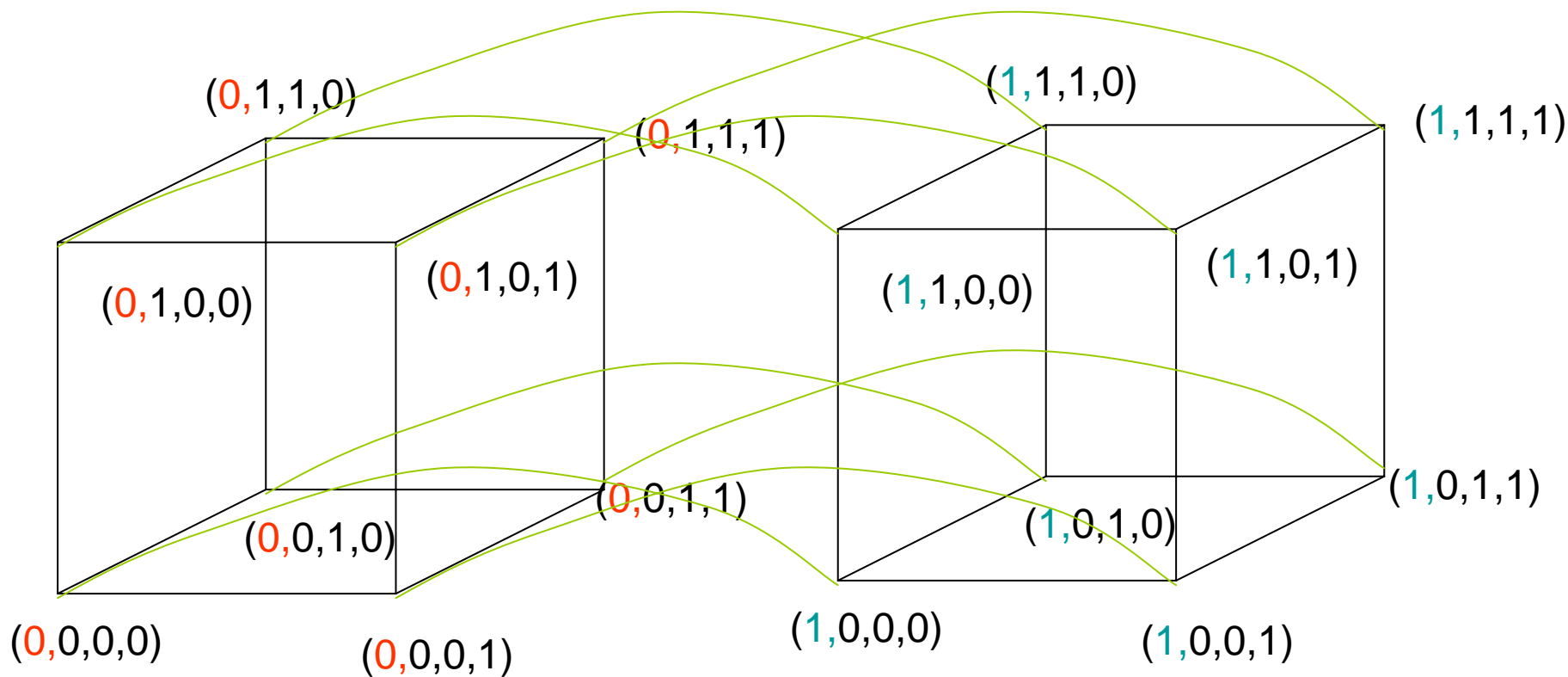


# 拡張は？



# 拡張

## 3キューブから4キューブへ



# 埋め込み能力は

- リング  
1次元の反射2進符号
- トーラス  
2次元の反射2進符号
- トリー  
1つのダミーノードを許すと可能

# 反射2進符号

符号間のハミング距離 = 1

通常2進    反射2進

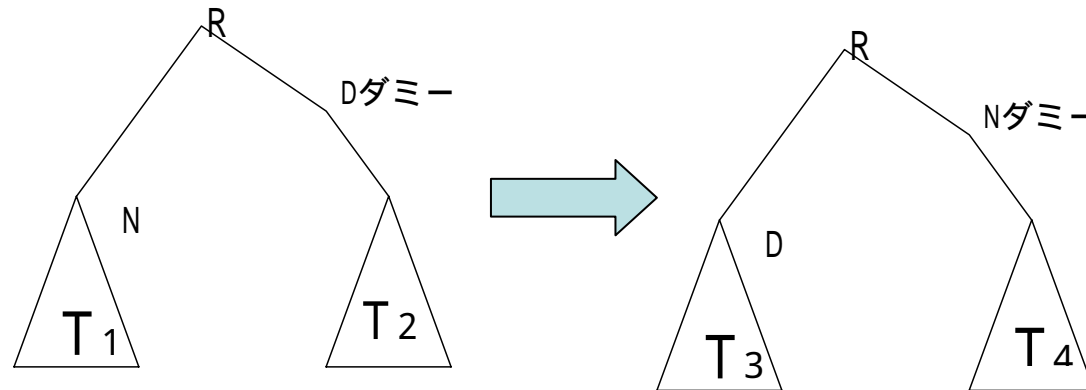
1 { 0 0 0  
0 0 1  
2 { 0 1 0  
1 { 0 1 1  
3 { 1 0 0  
1 { 1 0 1  
2 { 1 1 0  
1 { 1 1 1

1 { 0 0 0  
0 0 1  
1 { 0 1 1  
1 { 0 1 0  
1 { 1 1 0  
1 { 1 1 1  
1 { 1 0 1  
1 { 1 0 0

0 0 0  
0 0 1  
-----  
0 1 1  
0 1 0  
-----  
1 1 0  
1 1 1  
1 0 1  
1 0 0

# トリーの埋め込み

証明

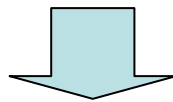


$n=k$

$R : (a_n a_{n-1} \dots a_N \dots a_D \dots a_1)$

$D : (a_n a_{n-1} \dots a_N \dots \bar{a}_D \dots a_1)$  で始まるダミー + T2

$N : (a_n a_{n-1} \dots \bar{a}_N \dots a_D \dots a_1)$  で始まる T1



$(a_n a_{n-1} \dots a_N \dots \bar{a}_D \dots a_1)$  で始まる T3

$(a_n a_{n-1} \dots \bar{a}_N \dots a_D \dots a_1)$  で始まるダミー + T4

が存在する。

$$R : ( a_n a_{n-1} \dots a_N \dots a_D \dots a_1 )$$

$$R' : ( a_n a_{n-1} \dots a_D \dots a_N \dots a_1 )$$

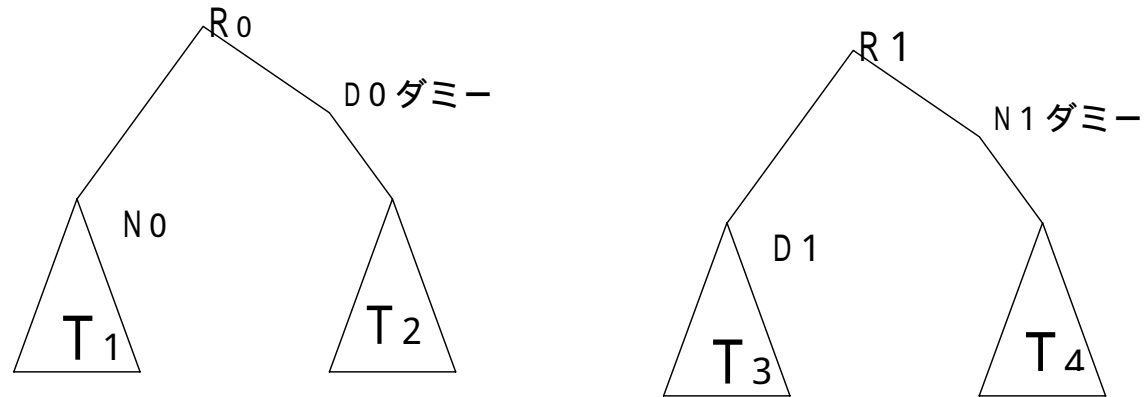
RとR'は1対1対応

$( a_n a_{n-1} \dots \bar{a}_N \dots a_D \dots a_1 ), ( a_n a_{n-1} \dots a_N \dots \bar{a}_D \dots a_1 )$ から始まる  
 トリーに重なりがなければ

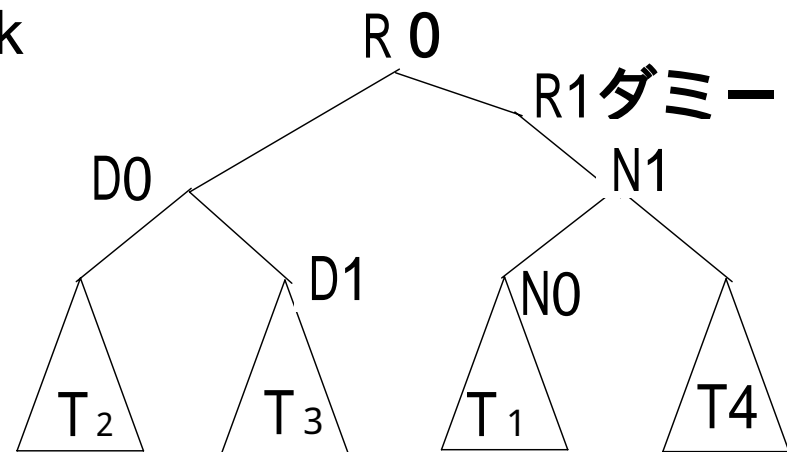
$( a_n a_{n-1} \dots \bar{a}_D \dots a_N \dots a_1 ), ( a_n a_{n-1} \dots a_D \dots \bar{a}_N \dots a_1 )$ から始まる  
 トリーに重なりがない。

# トリーの埋め込み

証明



$n=k$



$n=k+1$

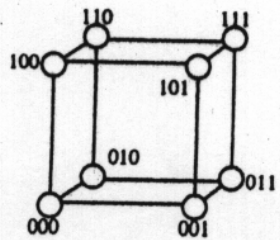


ベース  $m \times n$  - キューブ

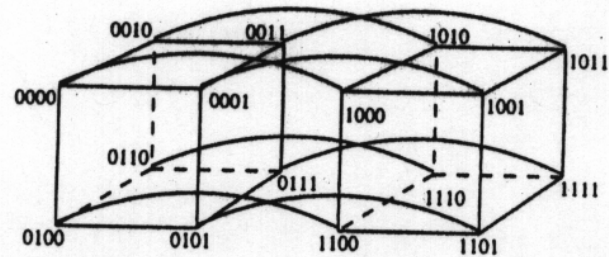
環立方体網 (CCC, Cube Connected Cycle)

$N = n \cdot 2^n$ 、次数は 3、直径は  $3n/2$

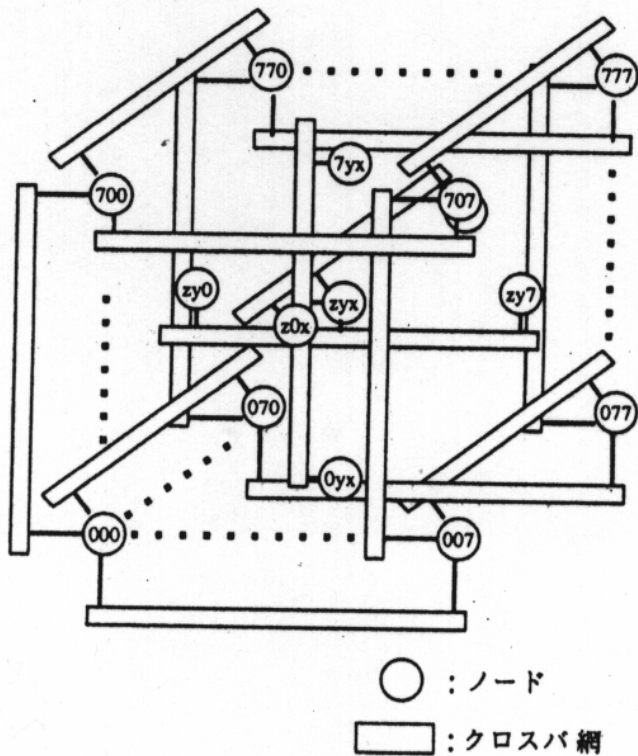
CCTcube網



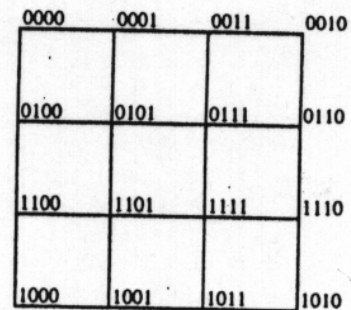
(a) 2進 3-キューブ



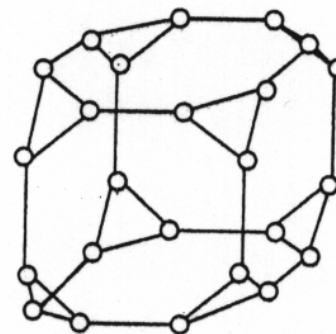
(b) 2進 4-キューブの合成



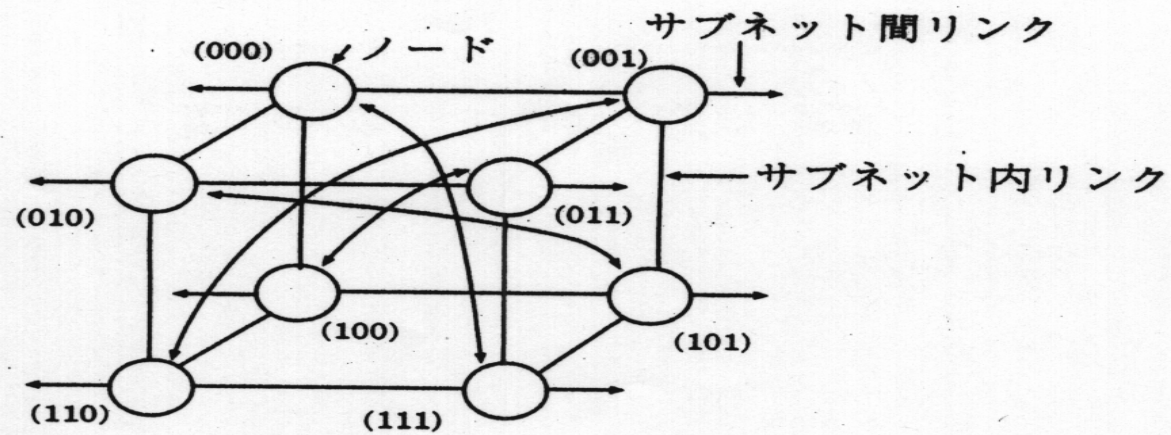
(d) ベース mn-キューブ



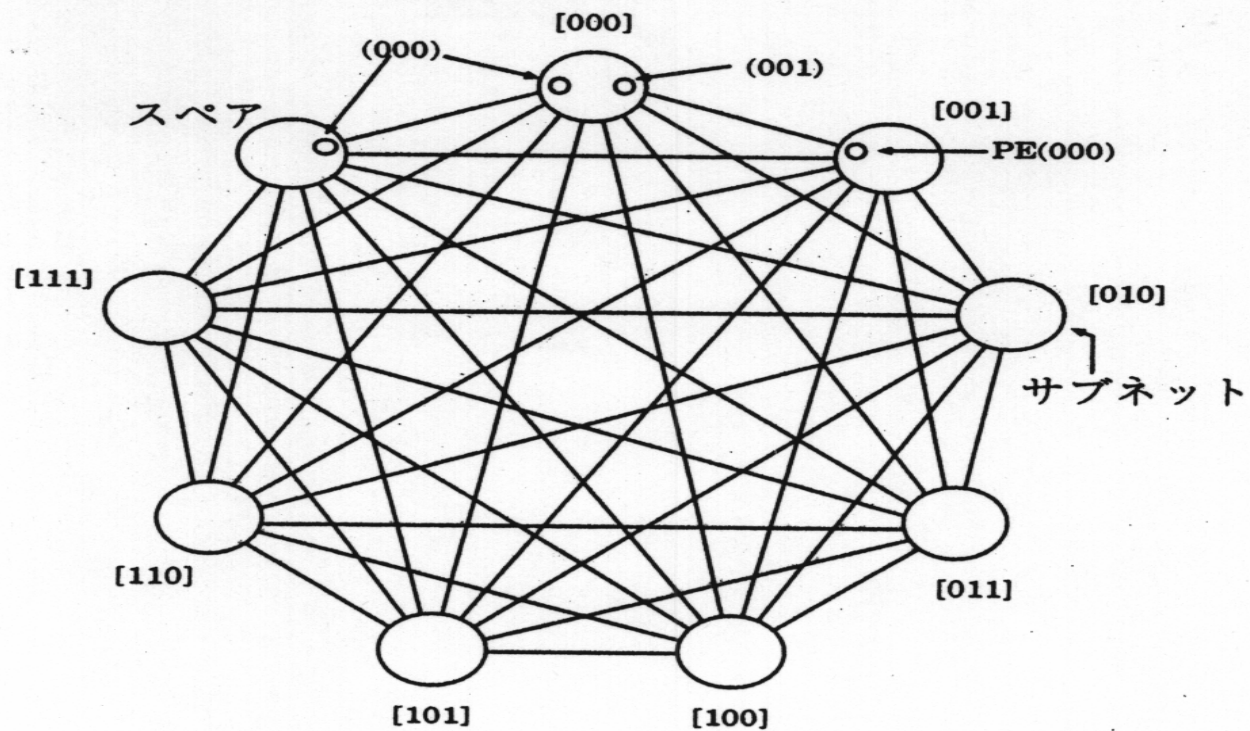
(c) 2進 4-キューブによるトーラスの実現



(e) 環立方体



(a) サブネット



(b) サブネット間の結合

図 2.8 CCT cube の構成

# ハイパクロス網

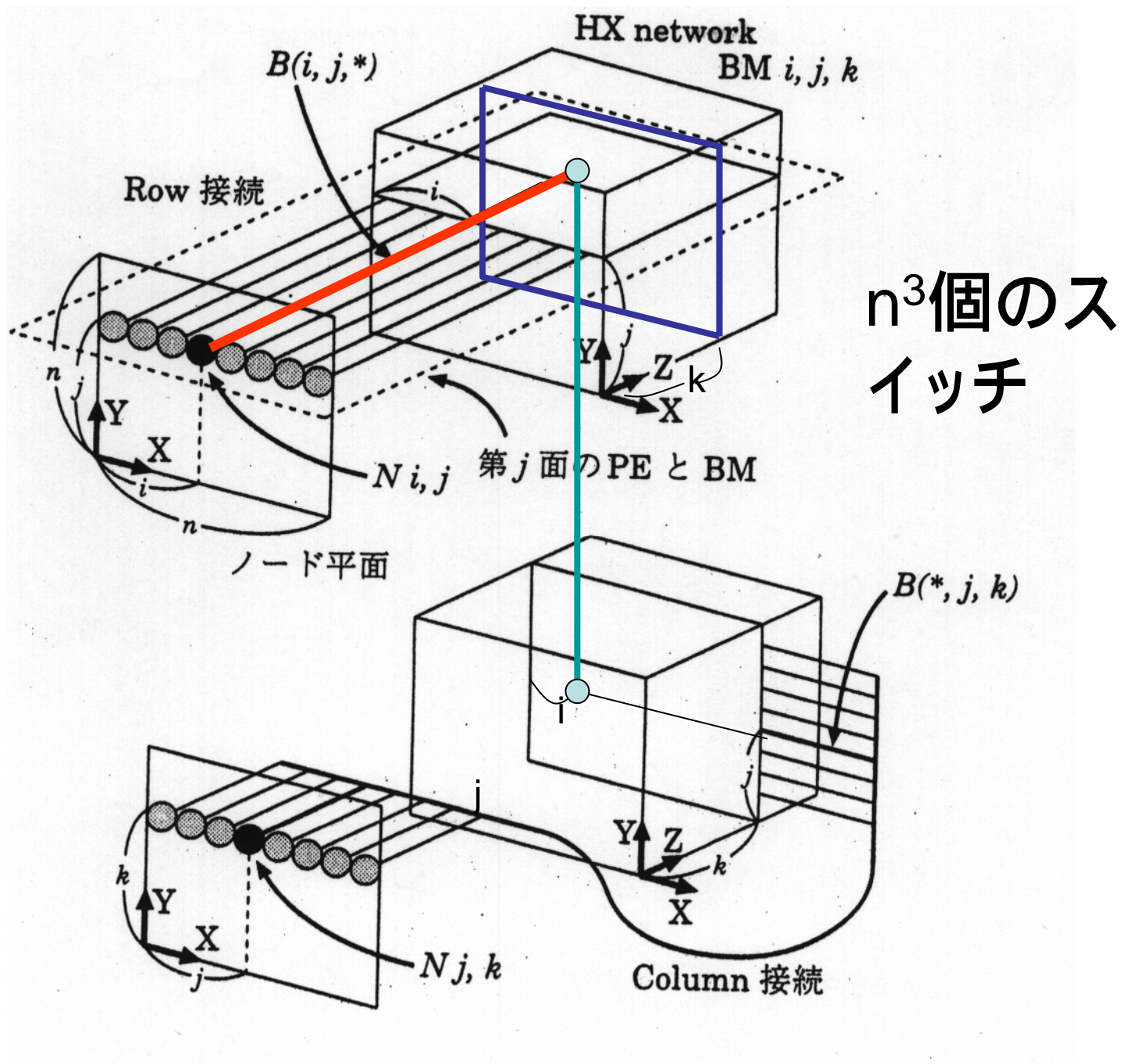
$$N=n^2$$

ノード番号を $N_{i,j}$

$N_{i,j}$ と $N_{j,k}$ および $N_{j,k}$ と $N_{i,j}$ が直接通信

任意のノード間通信 $N_{i,j}$ と $N_{k,l}$ の通信

$N_{i,j}$   $N_{j,k}$   $N_{k,l}$  直径は2



- ・ リング網の埋込み

$N_{k,j} (j=1,2,\dots,k)$ 、 $N_{i,k} (i=1,2,\dots,k)$  の

$2k-1$ 個のノード

$(1,k)$ 、 $(k,k)$ 、 $(k-1,k)$ 、 $(k,k-1)$ 、

$(k-2,k)$ 、 $(k,k-2)$ 、 $(k-3,k)$ 、 $\dots$ 、 $(k,1)$

とつなぎ、

$k+1$ について

$(1,k+1)$ 、 $(k+1,k+1)$ 、 $(k,k+1)$ 、 $(k+1,k)$ 、

$(k-1,k+1)$ 、 $(k+1,k-1)$ 、 $(k-2,k+1)$ 、 $\dots$ 、

$(k+1,1)$

とつなぐ。

つなぎ目の  $(k,1)$  と  $(1,k+1)$  は結合可

- ・ 2次元トラスの埋込み

2次元トラスのノード番号： $T_{x,y}$

$T_{x,y} \cdots N_{x,y}$ ： $x+y$ が偶数のとき

$T_{x,y} \cdots N_{y,x}$ ： $x+y$ が奇数のとき

T平面で $x+y$ が奇数：

隣接する4つのノードでの  
値（ $X+Y$ 値）はすべて偶数

$T_{x,y}$ のみが $N_{y,x}$ に写像。

N平面ではすべてHX結合条件満足

T平面で $x+y$ が偶数：

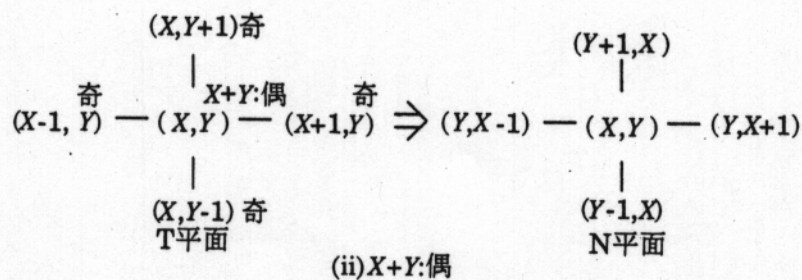
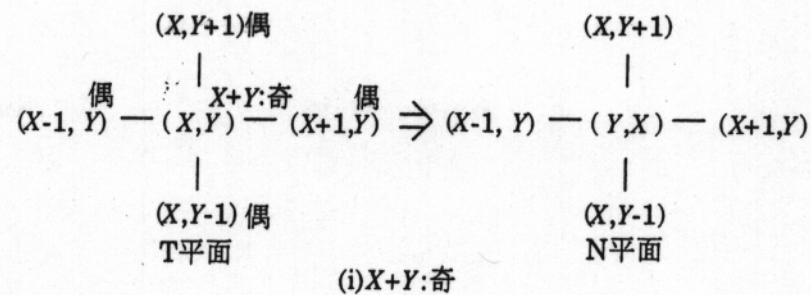
$T_{x,y}$ 以外すべて奇数。



$k=4$	16 4,1	14 4,2	12 4,3	11 4,4
$k=3$	9 3,1	7 3,2	6 3,3	13 3,4
$k=2$	4 2,1	3 2,2	8 2,3	15 2,4
$k=1$	1 1,1	2 1,2	5 1,3	10 1,4
	$k=1$	$k=2$	$k=3$	$k=4$

右上の番号: リング網のノード番号

(a) HX網へのリングの埋込み



4,1	4,2	4,3	4,4
3,1	3,2	3,3	3,4
2,1	2,2	2,3	2,4
1,1	1,2	1,3	1,4

T平面

1,4	4,2	3,4	4,4
3,1	2,3	3,3	4,3
1,2	2,2	3,2	2,4
1,1	2,1	1,3	4,1

N平面

(iii) T平面のN平面への写像

(b) HX網への2次元トーラスの埋込み



- ・ ハイパーキューブの埋込み

$$N=2^{2^k} (n=2^k)。$$

kビットの反射 2 進符号

ある数*i*とハミング距離 1 の数*j*

は*k*個存在

*j* - *i* : 奇数

# 反射2進符号

符号間のハミング距離 = 1

通常2進    反射2進

1	{	0 0 0	1	{	0 0 0		0 0 0	1
		0 0 1			0 0 1		0 0 1	2
2	{	0 1 0	1	{	0 1 1		0 1 1	3
1	{	0 1 1	1	{	0 1 0		0 1 0	4
3	{	1 0 0	1	{	1 1 0		1 1 0	5
1	{	1 0 1	1	{	1 1 1		1 1 1	6
2	{	1 1 0	1	{	1 0 1		1 0 1	7
1	{	1 1 1	1	{	1 0 0		1 0 0	8

2kビットの数

$a(i)$ 、 $b(i)$

2kビット長の符号： $a(i)+b(j)$

ノード $H_{a(i)+b(j)}$ を

ハイパクロス網ノード $N_{x,y}$ に写像

$x=i$ 、 $y=j$ ： $i+j$  奇数時

$x=j$ 、 $y=i$ ： $i+j$  偶数時

$i+j$  が奇数のとき

(  $i$  が偶数 / 奇数、  $j$  が奇数 / 偶数 )

$H_{a(i)+b(j)}$  は  $N_{i,j}$  に写像

$H_{a(i)+b(j)}$  と結合しているのは、

$b(j)$  を固定して考えると

$k$  個の  $H_{a(r)+b(j)}$ 、 $r-i$  : 奇数

$a(i)$  を固定して考えると

$k$  個の  $H_{a(i)+b(s)}$ 、 $s-j$  : 奇数

(  $i$  : 偶、 $j$  : 奇、 $r$  : 奇、 $s$  : 偶 ) または

(  $i$  : 奇、 $j$  : 偶、 $r$  : 偶、 $s$  : 奇 ) )

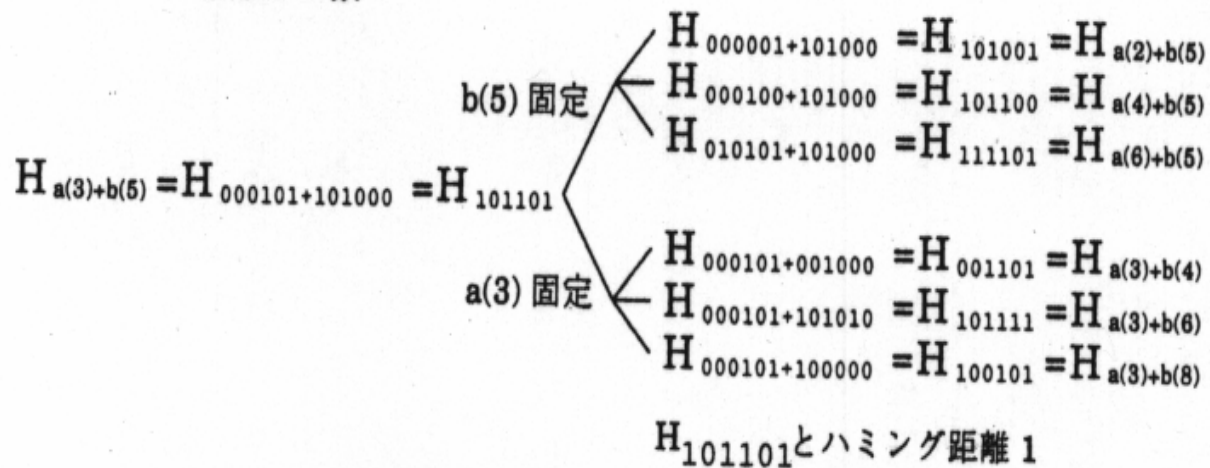
であり、 $r+j$ 、 $i+s$  は偶数

ノードは  $N_{j,r}$ 、 $N_{s,i}$  に写像 :  $N_{i,j}$  と結合

数	反射 2 進符号	$a(i)$	$b(i)$
1	000	000000	000000
2	001	000001	000010
3	011	000101	001010
4	010	000100	001000
5	110	010100	101000
6	111	010101	101010
7	101	010001	100010
8	100	010000	100000

→ 数 3 とハミング  
 距離 1 の数  
 ..... 数 5 とハミング  
 距離 1 の数

0 挿入      0 挿入



(c) HX 網へのハイパキューブの埋込み

図 2.10 HX 網への各種網の埋込み

## ( 7 ) シャフル網系

シャフル・エクスチェンジ網

de Bruijn 網

次数：4、直径： $\text{Log}_2 N$

ハミルトンパス：

000,001,011,111,110,101,010,100

## リング網の埋込み

n ビット列： $X_i = (x_{i1}, x_{i2}, \dots, x_{in})$

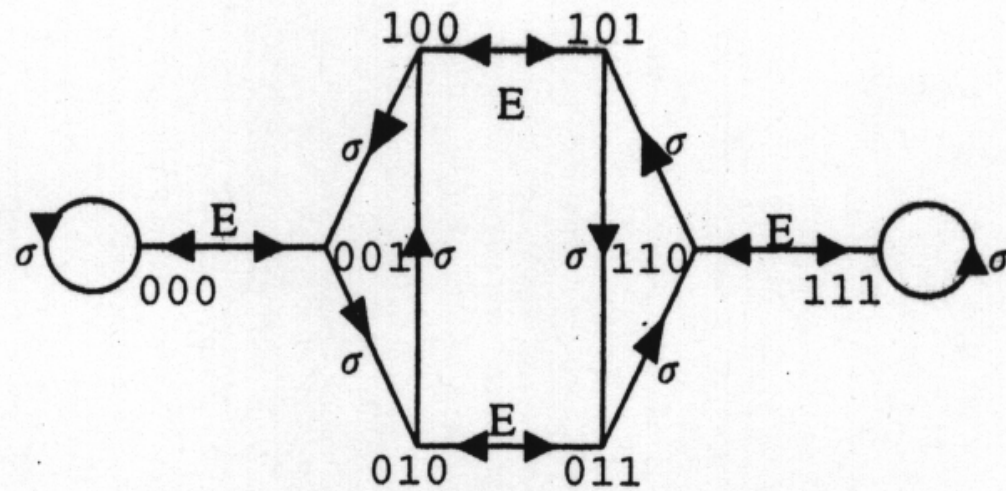
$$X_i = EX_i$$

$$EX_i = X_i$$

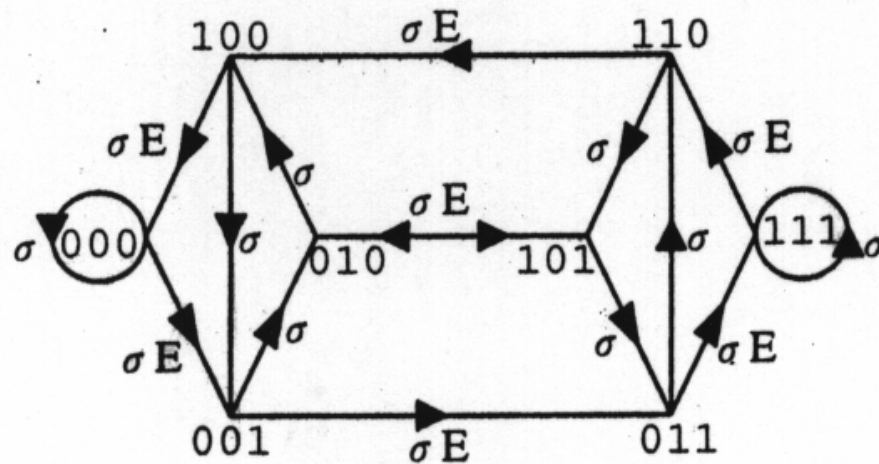
$X_i$  の共役 n ビット列： $X_i = (\sim x_{i1}, x_{i2}, \dots, x_{in})$

$(X_1, X_2, \dots, X_k)$  : サイクル

$$X_1 = X_2, \quad X_2 = X_3, \dots, \quad X_{k-1} = X_k, \quad X_k = X_1$$



(a) シャフル・エクスチェンジ網



(b) de Bruijn網

相互隣接サイクル  $C_1, C_2$

$$C_1 = (X_1^1, X_2^1, \dots, X_k^1), C_2 = (X_1^{1'}, X_2^2, \dots, X_m^2)$$

$C_1, C_2$ を1つの拡張サイクルに合成可能

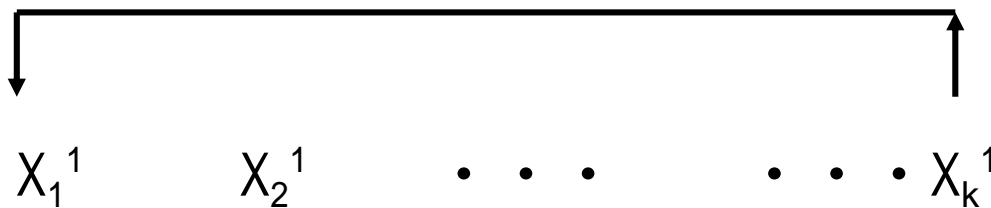
共役

拡張サイクル: E使用可能

$$X_1^1 = EX_1^1 = X_2^1$$

$$EX_1^1 = X_1^1 = X_2^2$$

$C_1$ サイクル



$C_2$ サイクル



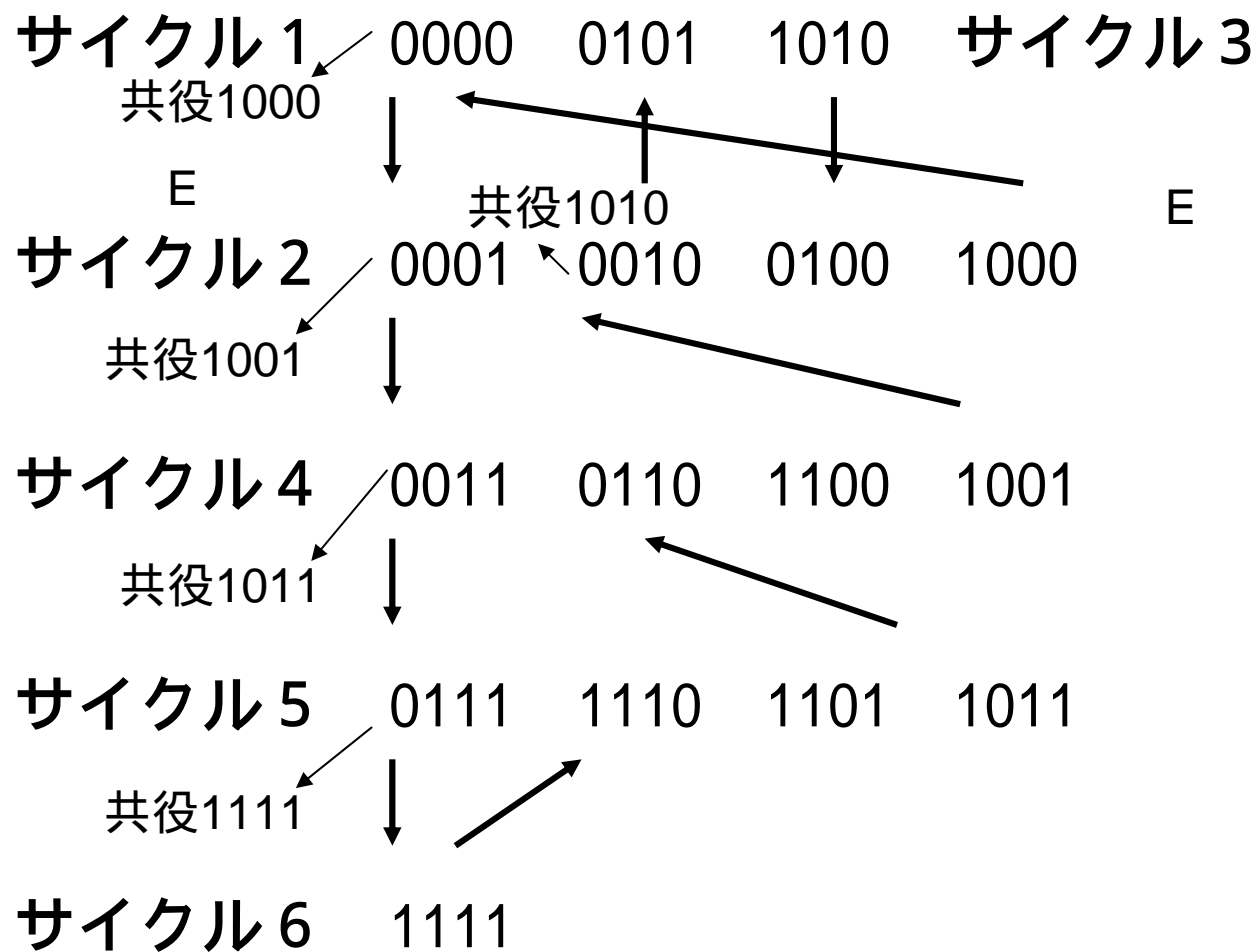
**重み**  $W(X_i) = \sum_{j=1}^n x_{ij}$

## アルゴリズム

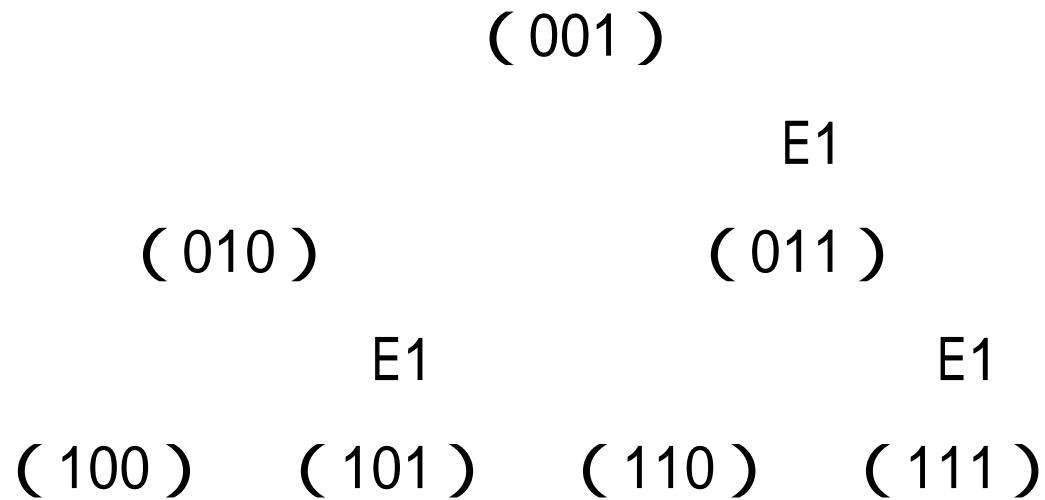
- (1) 重み 0 の  $n$  ビット列  $(0, 0, \dots, 0)$  からスタート
- (2) 重み  $K-1$  以下のサイクルをマージした拡張サイクル  $EC^{K-1}$  の要素  $X$  と共役な  $X$  を要素に持つ重み  $K$  のサイクルを選びマージし,  $EC^K$  とする.
- (3)  $K=n$  まで (2) を繰り返す.

例

重み 0	サイクル 1	0000				
重み 1	サイクル 2	0001	0010	0100	1000	
重み 2	サイクル 3	1010	0101			
	サイクル 4	0011	0110	1100	1001	
重み 3	サイクル 5	0111	1110	1101	1011	
重み 4	サイクル 6	1111				



## トリー網の埋込み



トーラス網：不可

de Bruijn網の一般形：

D進n桁  $(D_n, D_{n-1}, \dots, D_1)$  が

$\downarrow (D_{n-1}, D_{n-2}, \dots, D_1, *)$  と結合

循環オメガ網 (circular omega)

オメガ網のスイッチをノードで置き換えた網 (後述)

Star Graph網

ノード数  $N: n!$

記号列の最初の文字と他の文字を  
入れ換えたノードとリンク

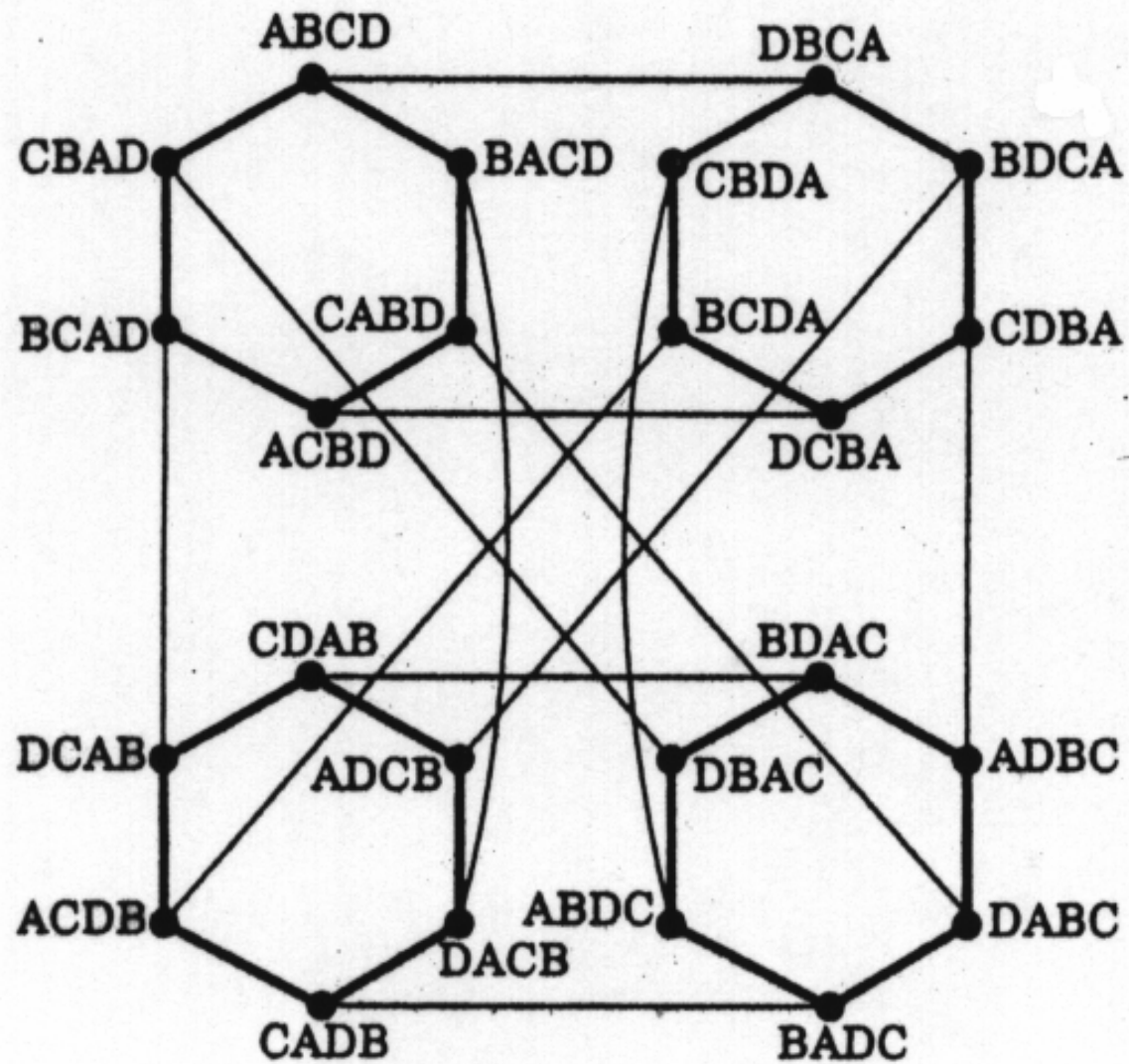
ノード ABCD : BACD、CBAD、DBCA と結合

次数 :  $n-1$

任意のノード番号から

ノード ABCDEFGH へのルーティング

記号 A, B, C, D, E, F, G, H の



(c)スターグラフ網の構成

ホームポジション：1,2,3,4,5,6,7,8

(たとえば、Aは1、Cは3、Hは8)

- ・与えられた送信ノードの最初の記号がAでない  
その記号に対応する送信ノード番号中のホーム  
ポジションにある記号と置き換えを行う。
- ・与えられた送信ノードの最初の記号がA  
適当な位置にある記号と入れ換える。
- ・以上の過程を繰り返す。

FDGBEHCA    HDGBEFCA    ADGBEFCH    DAGBEFCH  
BAGDEFCH    ABGDEFCH    GBABEFCH    CBADEFGH

ABCDEFGH

記号列のサイクル：サイクル内の入れ換えホームポジションへ

FDGBEHCA：

(F,H,A)と(D,B)と(G,C)サイクルの形成

ルーティング回数 $N_r$

$N_r = C + m - 2$ ：送信ノード番号の先頭がAでないとき

$C + m$ ：送信ノード番号の先頭がAのとき

FDGBEHCAの例： $C=3, m=7, N_r:8$

サイクルの決定

FDGBEHCAでは、サイクル1：(F,H,A)、



サイクル1 : (F,H,A)

サイクル2 : (D,B)、サイクル3として : (G,C)

サイクル1での置換

FDGBEHCA - >

HDGBEFCA - > ADGBEFCH

サイクル1 : サイクル長3、

2回の置換でADGBEFCH : 一般にもとのサイクル長-1

サイクル2に統合

ADGBEFCH - > DAGBEFCH - >

BAGDEFCH - > ABGDEFCH

Aをサイクル2に統合 : サイクル長3、

3回の置換でABGDEFCH : 一般にもとのサイクル長 + 1

サイクル3に統合 :

ABGDEFCH - > GBADEFCH - >

CBADEFCH - > ABCDEFGH

サイクル3、サイクル長3、

3回の置換でABCDEFGH : 一般にもとのサイクル長 + 1

先頭がAでないとき

置換回数： $m-1+C-1$ で $C+m-2$

先頭がAのとき： $C+m$

C-1個のサイクルで+1の置換

最初のサイクルで置換-1

直径

先頭：A

サイクル： $(D, B, C)$ 、 $(F, E)$ 、 $(H, G)$ ： $\lfloor (n-1)/2 \rfloor$

サイクル長：総数が $n-1$ のとき  
 $\lfloor 3(n-1)/2 \rfloor$

## **2 . 3 静的網のルーティング、 デッドロック**

### **2 . 3 . 1 パケット交換方式**

蓄積交換 (store and forward)

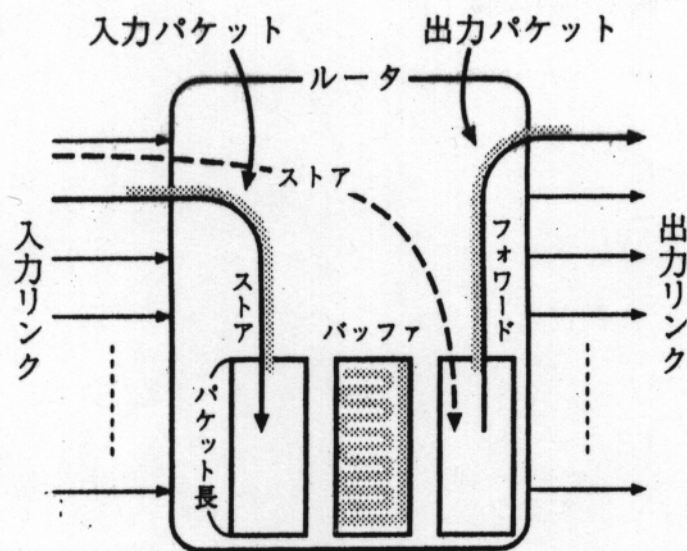
ワームホール

バーチャルカットスルー

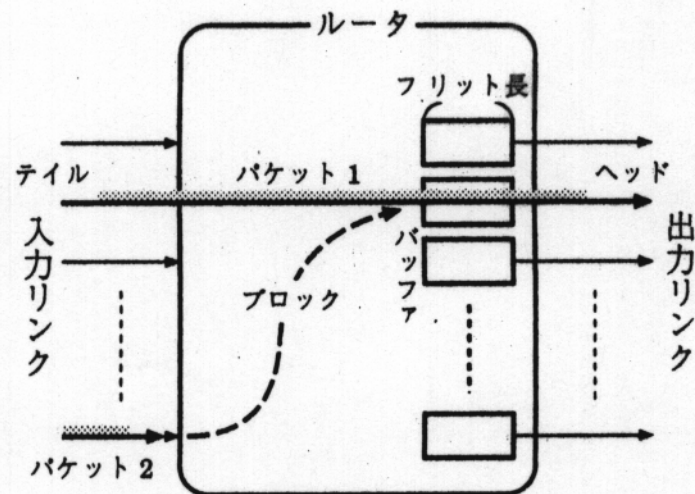
### **2 . 3 . 2 デッドロック一般論**

デッドロックの発生条件

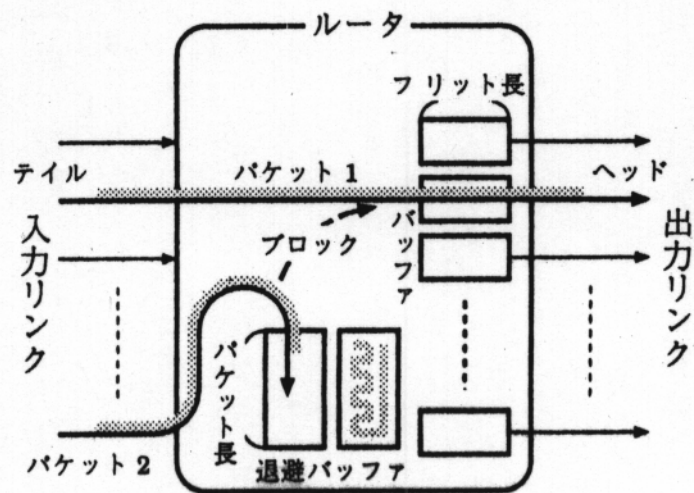
排他的利用 (Exclusive control)



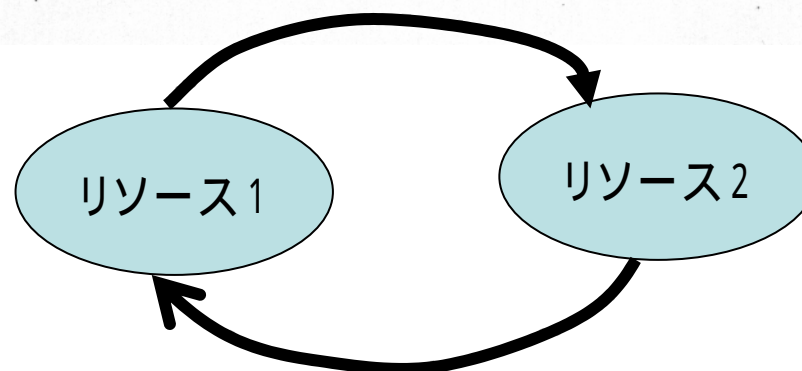
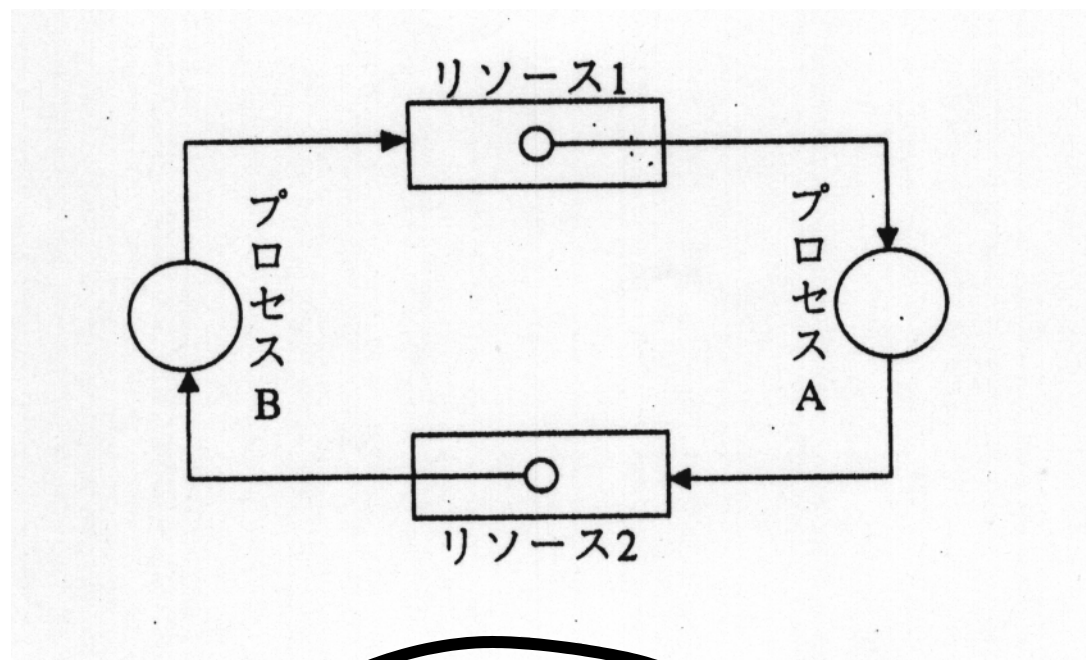
(a) 蓄積交換



(b) ワームホール



(c) バーチャルカットスルー



非妥協的リソース待ち (hold-wait)

横取り不可 (non-preemption)

リソース利用におけるサイクルの存在

デッドロック対処手段

( 1 ) デッドロック検出 (detection)

( 2 ) デッドロック回避 (avoidance)

銀行家のアルゴリズム

将来デッドロックが生じる可能性のある

リソース割付けの禁止

システム状態：安全 (safe) 状態に保つ

安全な状態にあるプロセス：

すべて実行終了可能

$P_0, P_1, P_2, \dots, P_m$  プロセス：

リソース利用最大値 ( $M_i$ ) が既知

システムが安全状態

現在リソース利用値： $U_i$

現在リソース要求残値  $B_i$ ： $(M_i - U_i)$

現在リソース要求値： $R_i$

空きリソース数  $T$ ：

(全リソース数 -  $\sum U_i$ )

$$T \geq R_i$$

	A 社 限度額 20 万円	銀行 資本 20 万円	B 社 限度額 15 万円
1 月	5 万円要求 受理	15 万円	
2 月		7 万円	8 万円要求 受理
3 月	4 万円要求 受理	3 万円	
4 月		2 万円	1 万円要求 受理
5 月	3 万円要求 拒否		
6 月		3 万円要求 拒否	
A, B 両社とも事業の推進不可能			



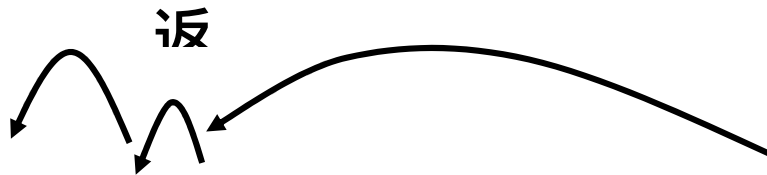
$$T \geq B_i$$

一方、 $T \geq B_i$  のとき受理される方式では次のようにデッドロックは生じない。

	A社	銀行	B社
	限度額 20 万円	資本 20 万円	限度額 15 万円
1 月	5 万円要求 受理	15 万円	
2 月		7 万円	8 万円要求 受理
3 月	4 万円要求 拒否 (将来の予想借入れ最大額 15 万円 > 銀行資金残高 7 万円)		
4 月		6 万円	1 万円要求 受理
5 月		3 万円	3 万円要求 受理
6 月		0 万円	3 万円要求 受理
7 月		15 万円	15 万円返済

## $P_i$ のリソース要求

T  $B_i$ のとき要求を受理



$P_1$	$P_2$	$P_3$	$P_4$	$P_5$	$P_6$	•	•	•	$P_8$	$P_9$	•	•	•	•
受拒	受拒	受拒	拒否	拒否	拒否	拒否	拒否	拒否	拒否	受	以後拒否	以後拒否	以後拒否	以後拒否
理否	理否	理否	理否	理否	理否	理否	理否	理否	理否	理	理	理	理	理

### ( 3 ) デッドロック防止 ( prevention )

リソースの共有化

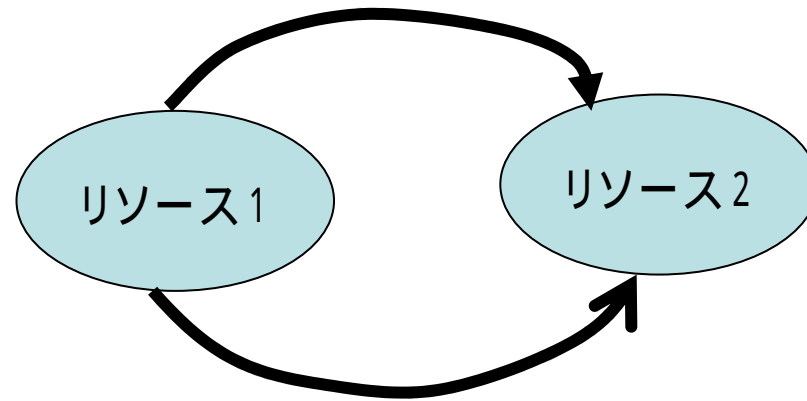
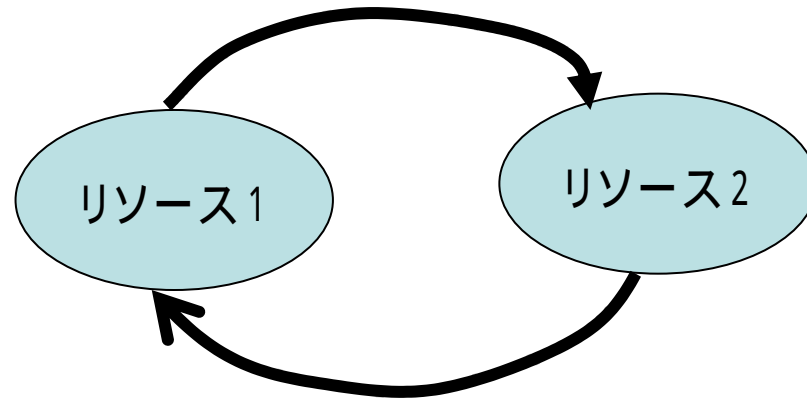
妥協的リソース待ち

- ・ 非受理時全リソース解放
- ・ 一括リソース要求

横取り可能

サイクルの防止

- ・ リソース数の増大
- ・ 代替リソースの提供
- ・ リソース番号順獲得



## 2.3.3 デッドロック

### (1) チャネル依存グラフ

チャネル依存グラフ  $D=G(C,E)$

$$E=\{(C_i,C_j) \mid R(C_i,n)=C_j\}$$

$R$  : ルーティング方式

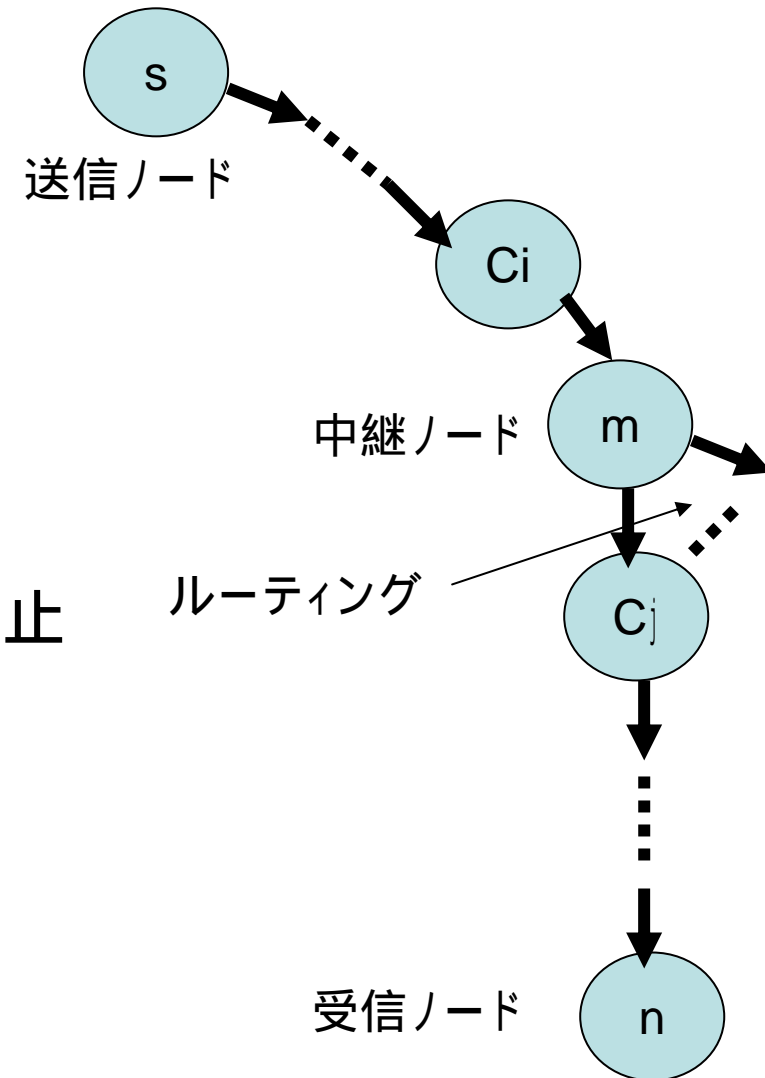
$n$  : 宛先  $n$

### (2) リング網でのデッドロック防止

チャネルの共有化

仮想化につながる考え

妥協的チャネル待ち



バーチャルカットスルーの考え方

退避バッファ：デッドロック問題

横取り可能チャネル

サイクルの防止

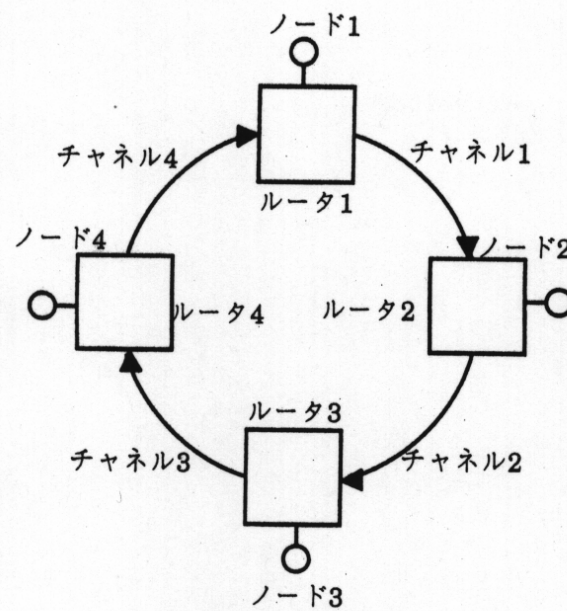
- ・ リングの切断
- ・ 仮想チャネル

仮想チャネル

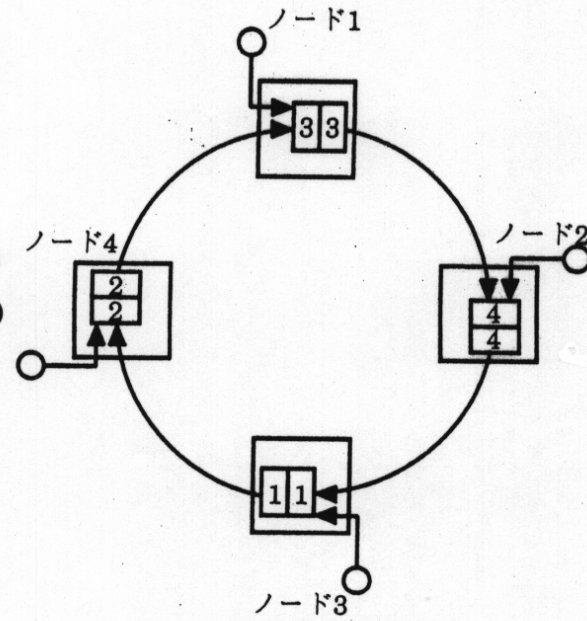
一意の番号付け

(チャネル番号： $C_1, C_2, \dots, C_8$ )

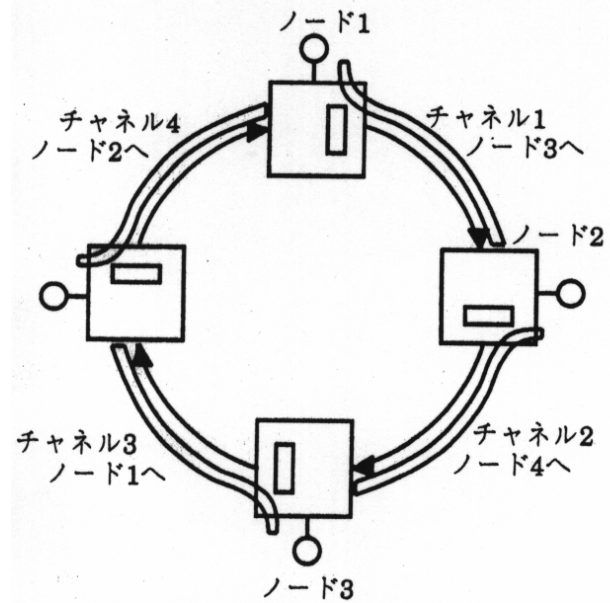
$C_1$ から $C_4$ ：下位リング



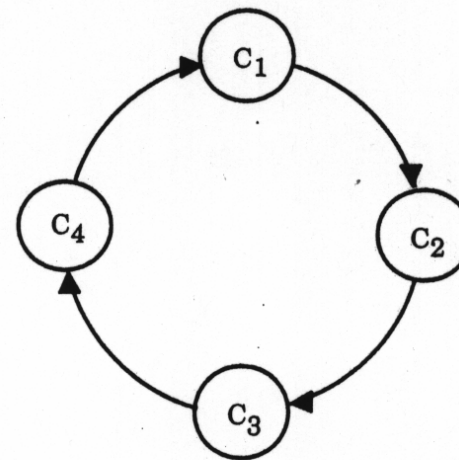
(a)リング網



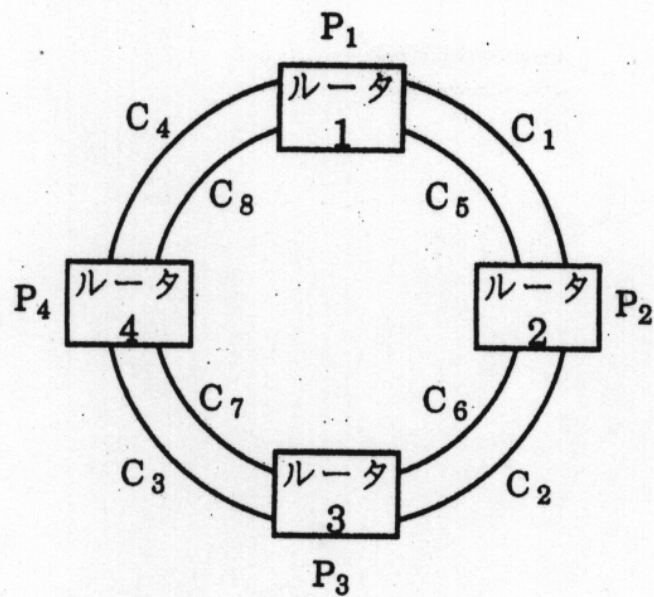
(b)蓄積交換方式



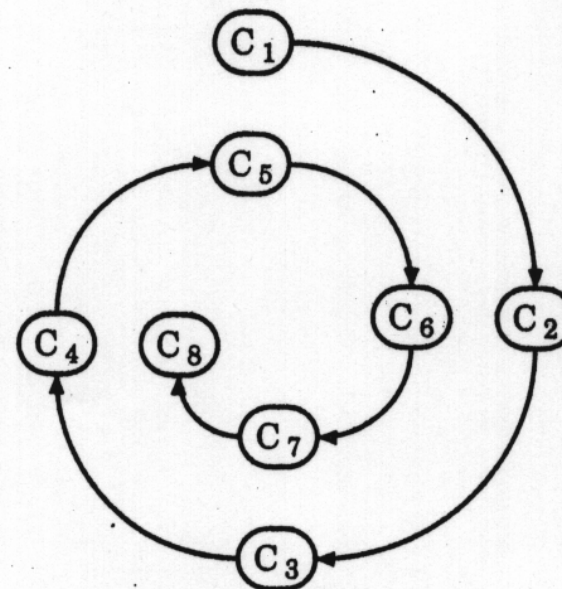
(c)ワームホール方式



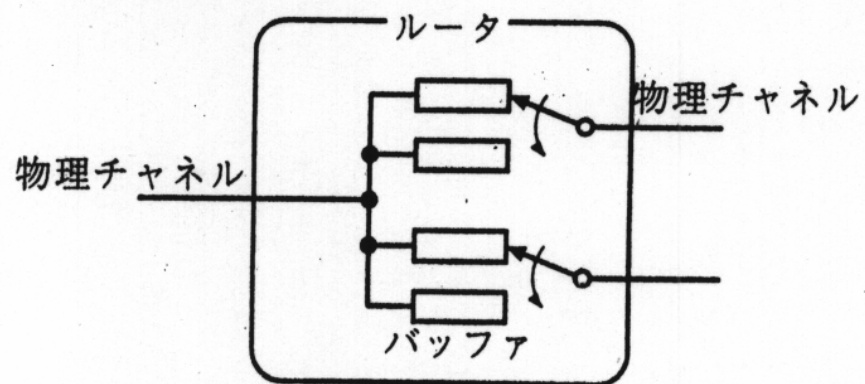
(d)チャンネル依存グラフ



(a) 仮想チャネル



(b) チャネル依存グラフ



(c) ルータの構成



$C_5$ から $C_8$ ：上位リング

ルーティング方式：チャネル番号が増大する方向

プロセッサ 1 から 3 への通信： $C_5$ と $C_6$

プロセッサ 4 から 2 への通信： $C_4$ と $C_5$

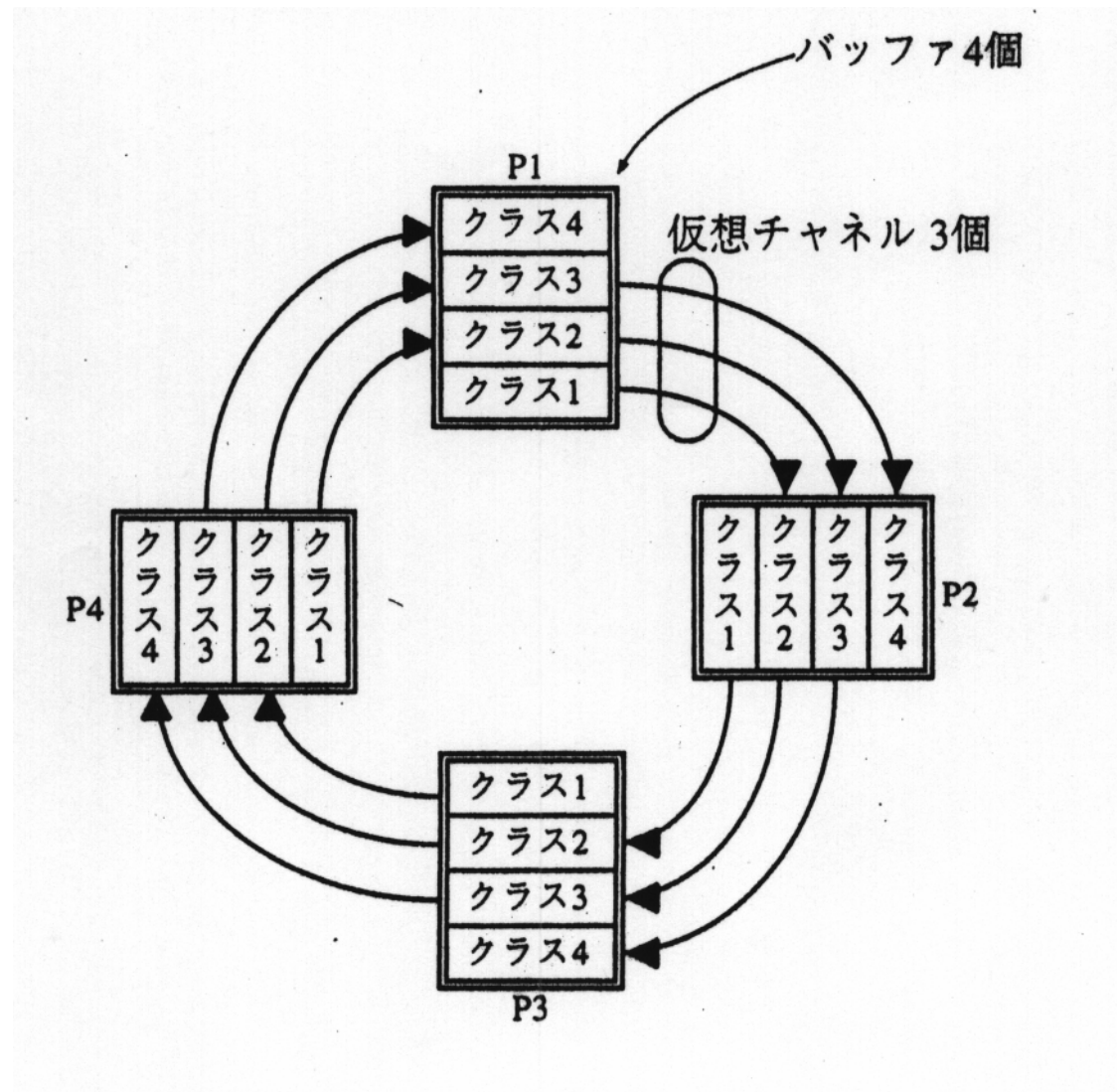
ルータ番号 $m$ に到着したパケットの受信先番号 $m$

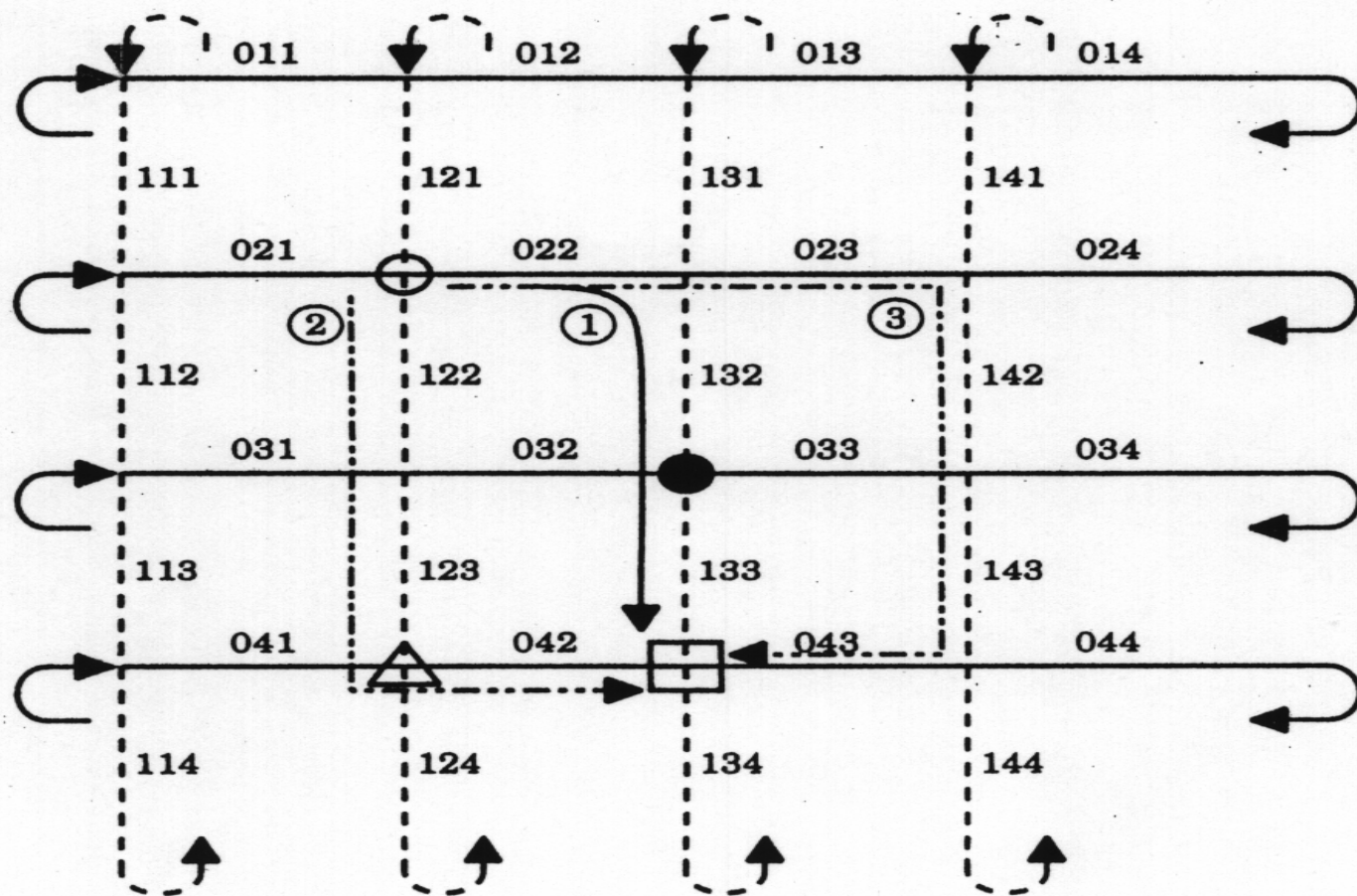
より小さいとき、下位リングを、

大きいときは上位リングを選択

( 3 ) トーラス網、ハイパキューブ網でのデッドロッ  
ク防止

次元順ルーティング





実線 : X方向

点線 : Y方向

① : 次元順ルーティング

② : 次元逆順ルーティング

③ : 迂回ルーティング

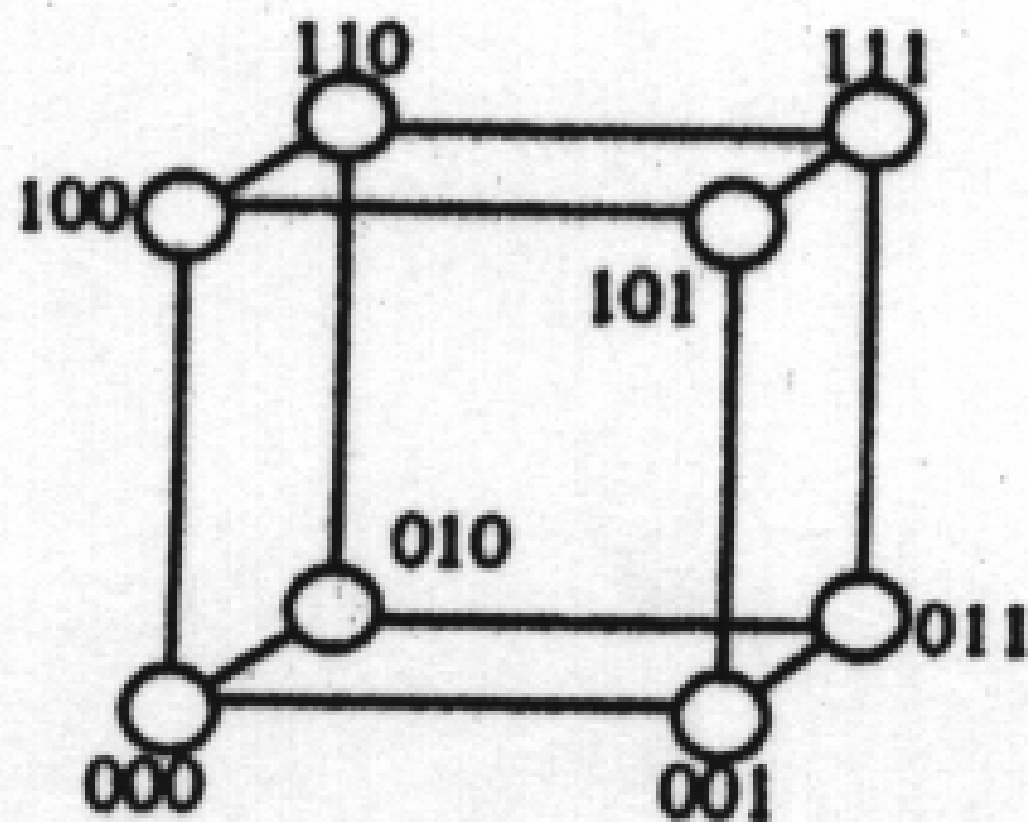
2次元トーラスの場合

ハイパーキューブの場合

次元順ルーティングでのデッドロック防止

すべてのチャンネルに一意的番号付けができ、

昇順に従ったチャンネル番号の獲得



(a) 2進 3-キューブ

## 2.3.4動的ルーティング

### 次元順ルーティング

デッドロック防止が可能

ルーティングが固定方式

トラフィックの偏り、耐故障に弱い

デッドロックを生じない方式

受信ノードに少なくとも辿り着く方式

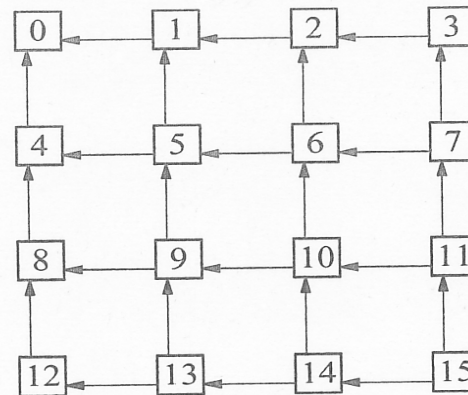
# (1) 仮想チャネルを使用した方法

(a) シンプルな方法

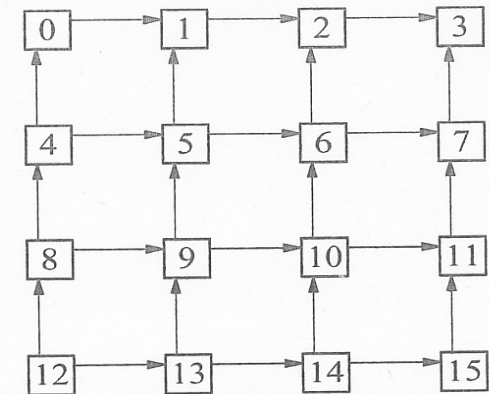
0 15 のとき: X+Y- ネット

1 0 のとき: X-Y- ネット

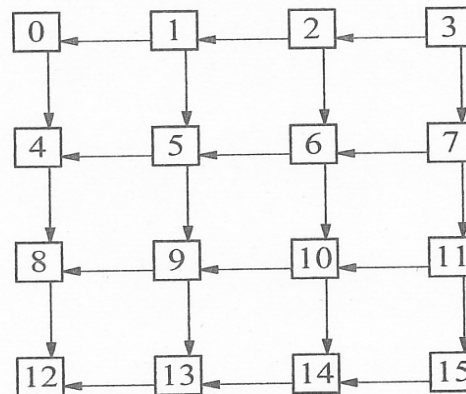
X-Y+ Virtual Network



X+Y+ Virtual Network



X-Y- Virtual Network



X+Y- Virtual Network

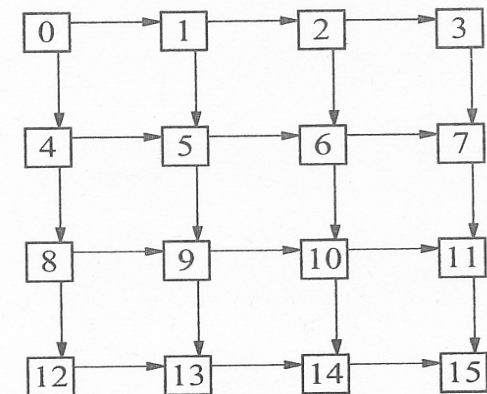


Figure 4.16. Virtual networks for a 2-D mesh.

( b ) 逆次元順ルーティング

ルータ間のチャネル：

仮想チャネルによって  $r$  多重

$r$  個のサブネットワークを構成：

クラス  $i$  サブネットワーク

パケットにクラスフィールド  $C$ ：初期値を 0

パケット：クラス番号  $C$  に対応した

サブネットワーク内で任意の次元方向



( X 方向、 Y 方向 ) にルーティング

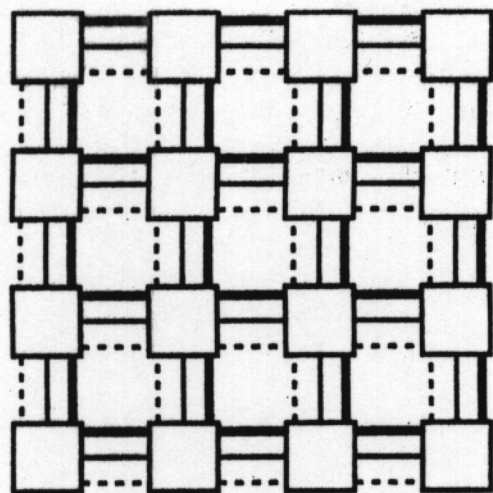
高位次元から低位次元へルーティングの切換え

が生じたとき( Y 方向後、X 方向に切り替え )

C 値を + 1

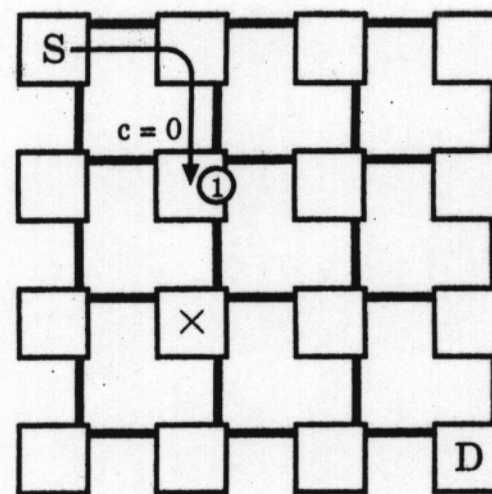
$C=r$  のとき以後はクラス  $r$  のサブネット

ワークで次元順

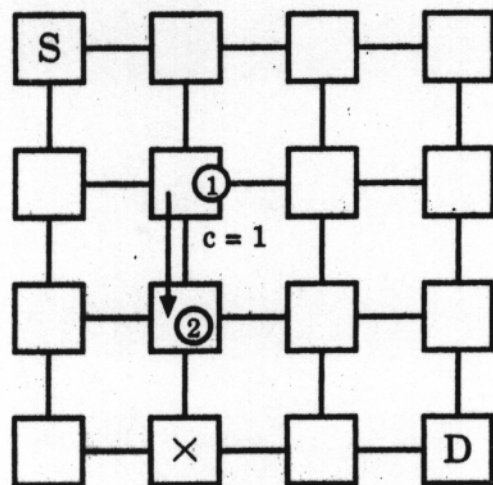


— : バーチャルチャネル 0  
 - - : バーチャルチャネル 1  
 . . . : バーチャルチャネル 2

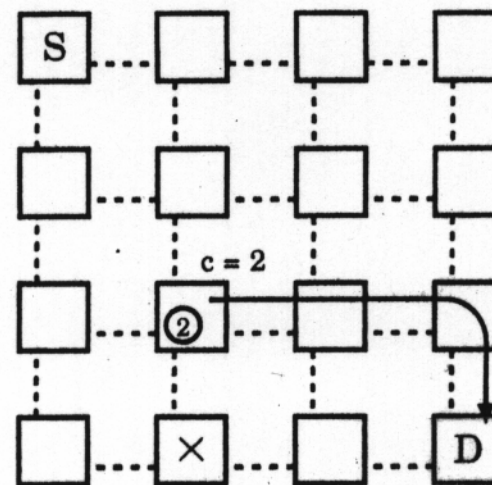
(a) バーチャルチャネル



(b) クラス 0 サブネットワーク



(c) クラス 1 サブネットワーク



(d) クラス 2 サブネットワーク

## ( 2 ) BlueGene/Lの方式

Escape Ring: 単方向、次元順ルーティング,

外部からQueueに2つ空きの時: 注入可能

デッドロックなし

内部からQueueに1つの空きの時: 投入可能

Adaptive Ring: Queueが一杯の時, Escape Ringに注入

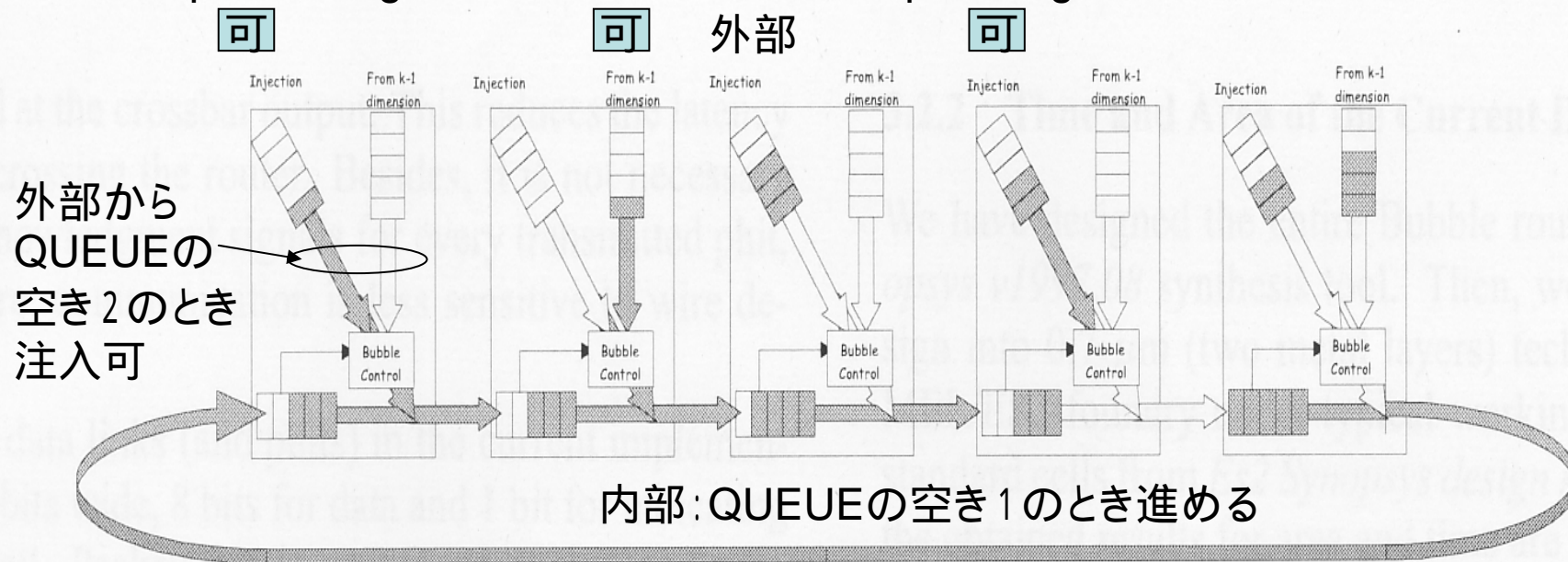
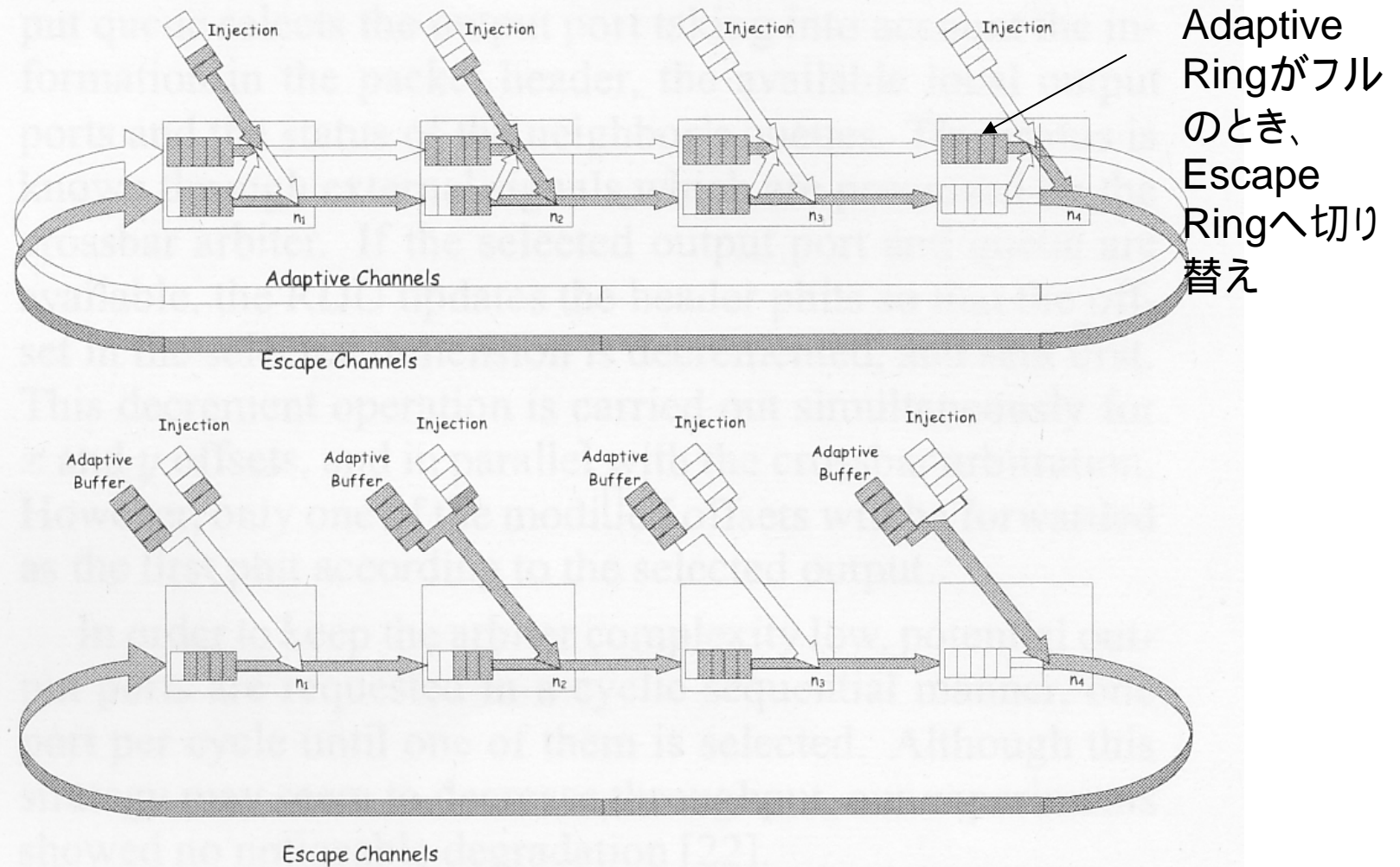


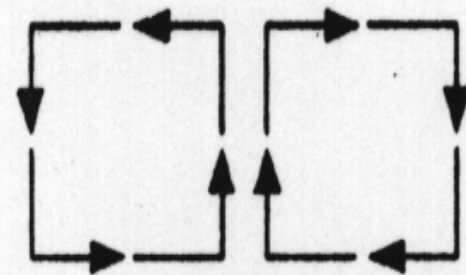
Figure 1. Deadlock avoidance in a ring by using Bubble flow control.



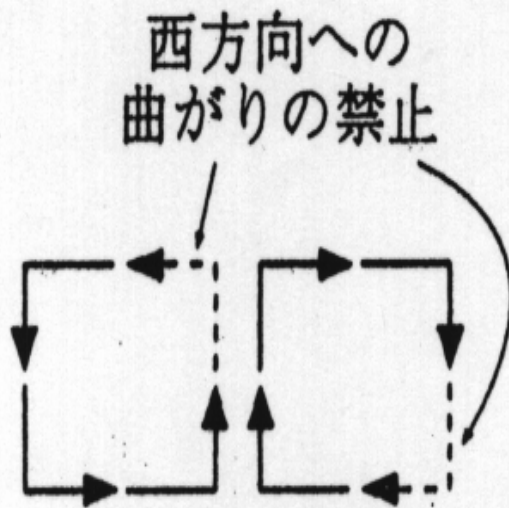
**Figure 2. Representation of adaptive and escape queues when all adaptive ones are busy and equivalent network.**

#### (4) 曲がり方向制限法 (turn model)

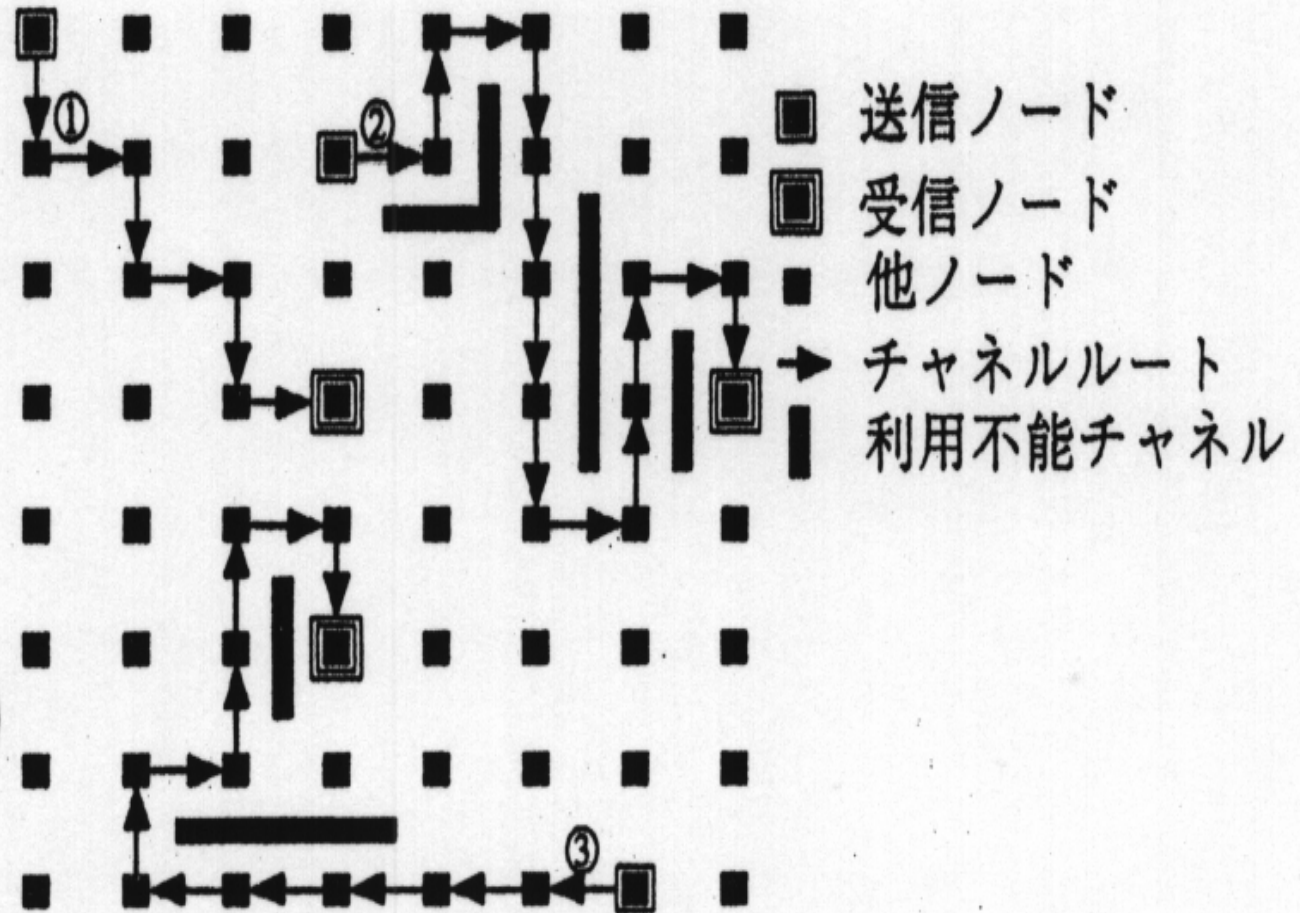
##### 「西」への曲がりを禁止するルーティング法



(a) サイクル形成



(b): サイクル除去



(c) ルーティング例

## **2 . 4 動的網**

### **2 . 4 . 1 クロスバ網**

集中制御方式

分散制御方式

### **2 . 4 . 2 多段結合網**

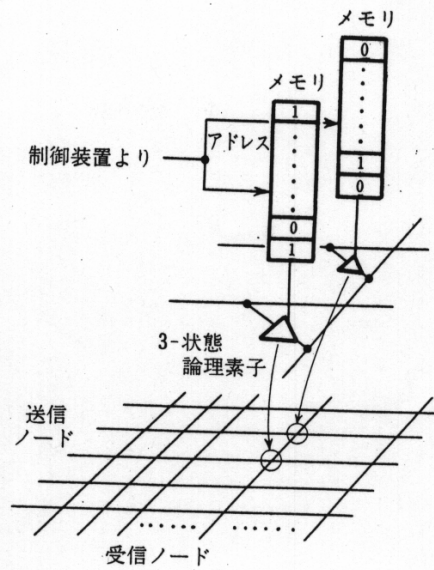
( 1 ) 完全非閉塞網

通信パターン数： $N!$

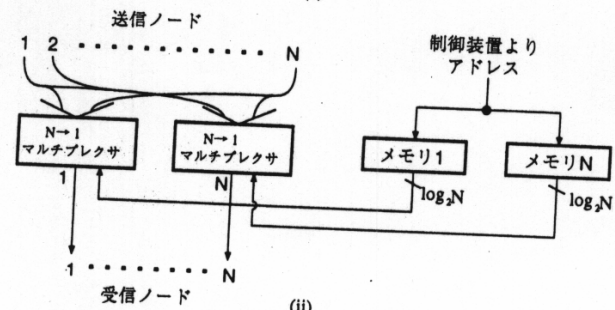
通信変更：ローカルにスイッチ変更で対処可

クロスバ網



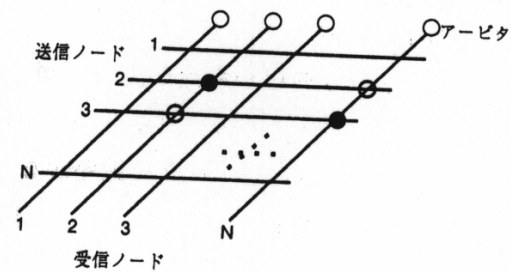


(i)



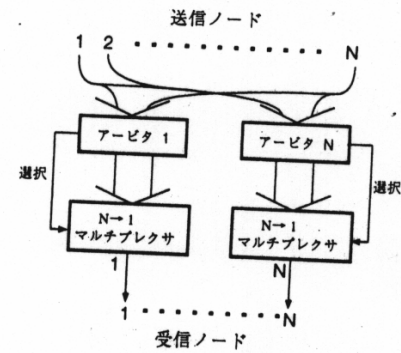
(ii)

(a) 集中制御方式例



(iii)

(b) 分散制御方式例



(iv)

図 2.20 クロスバ網

### 3 ステージ Clos 網

m  $2n-1$  のとき非閉塞：存在定理

### スイッチ数

$$N_s = (2n-1) \{ (N/n)^2 + 2N \}$$

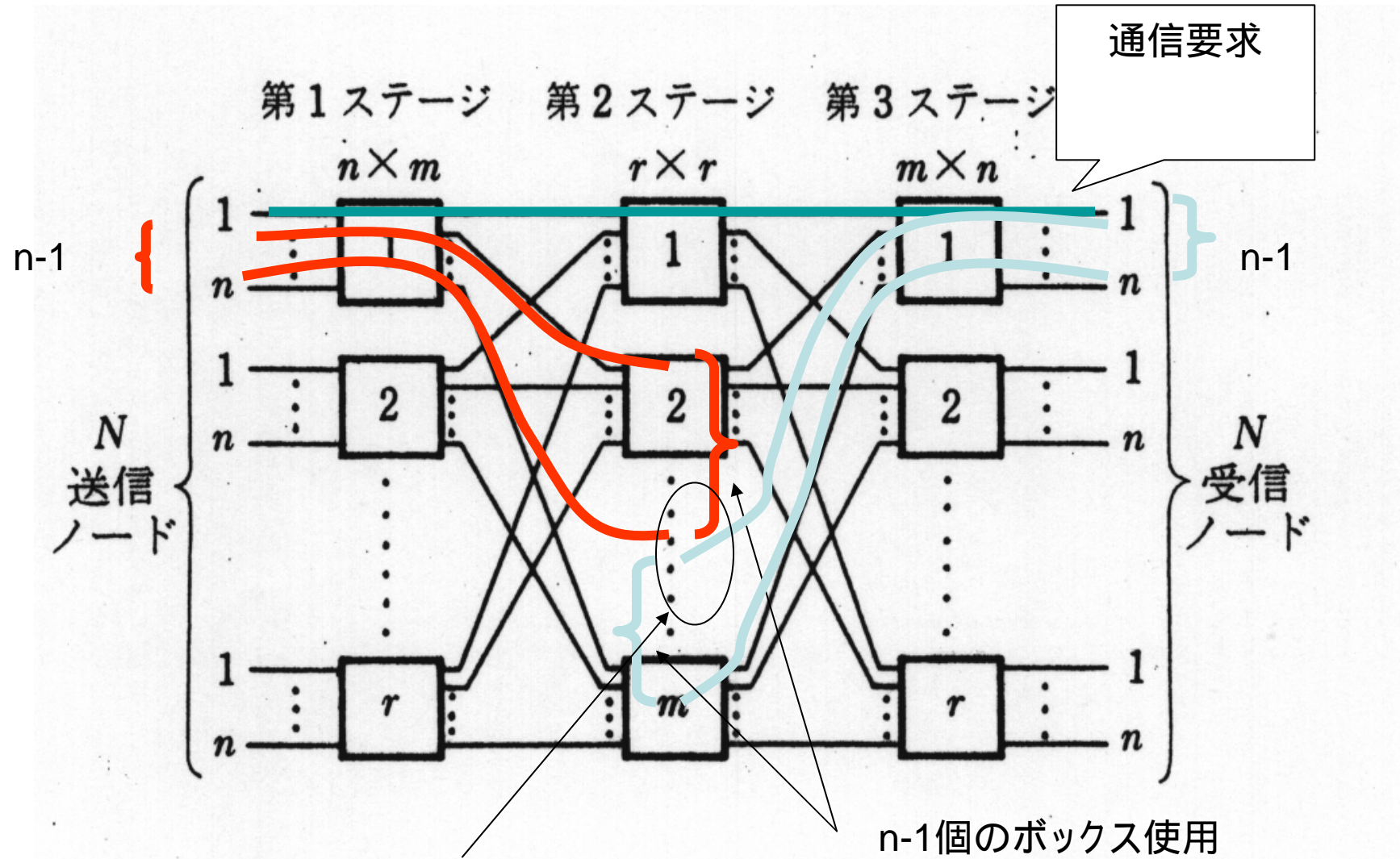
$$\text{Min } N_s = 4\sqrt{2}N^{3/2} - 4N$$

$N = 256$  の時

22146個      クロスバ：65536個

### ルーティング困難





## ( 2 ) 再構成型非閉塞網

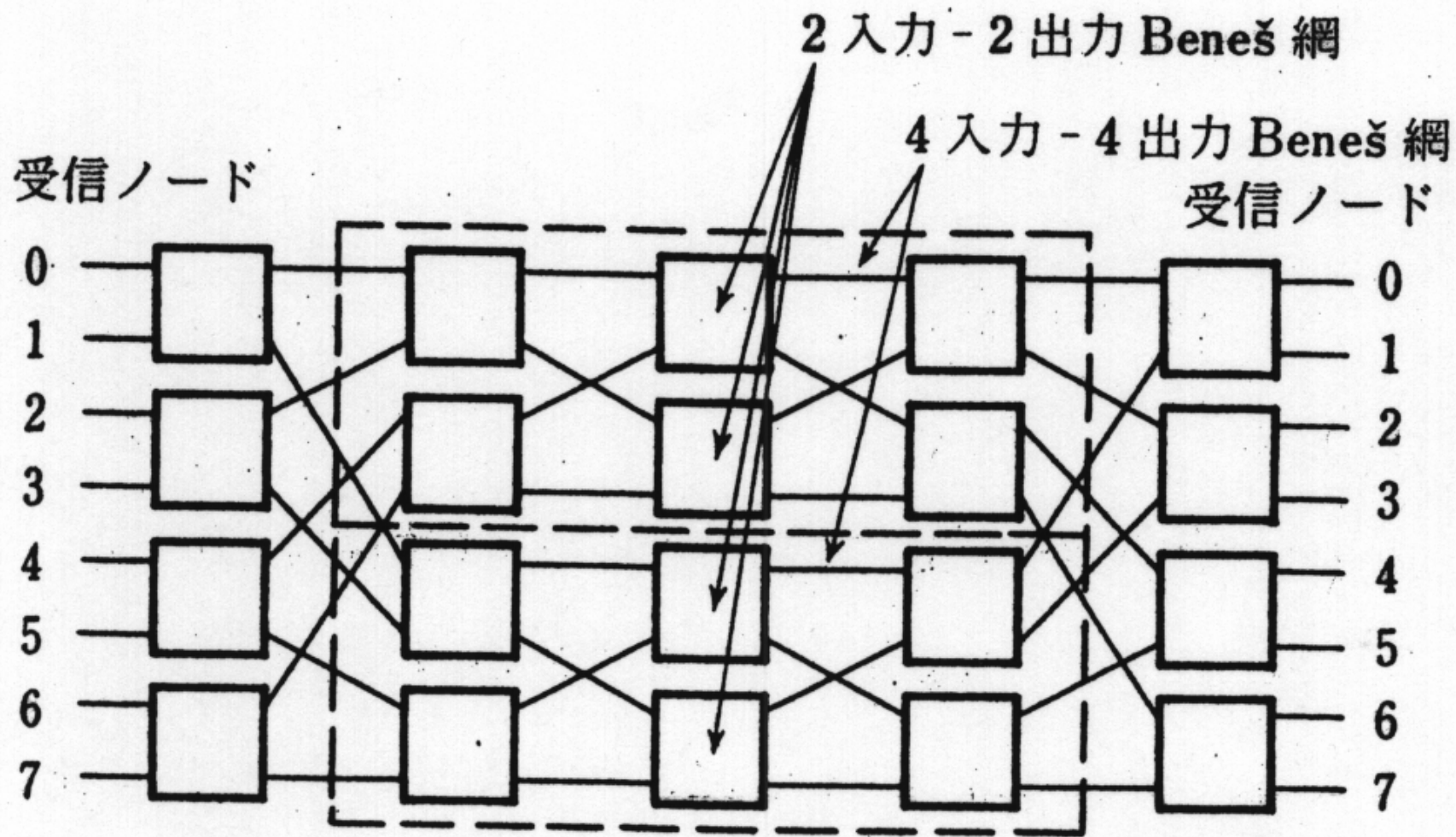
通信パターン数 :  $N!$

通信変更 : ローカルにスイッチ変更で対処不可

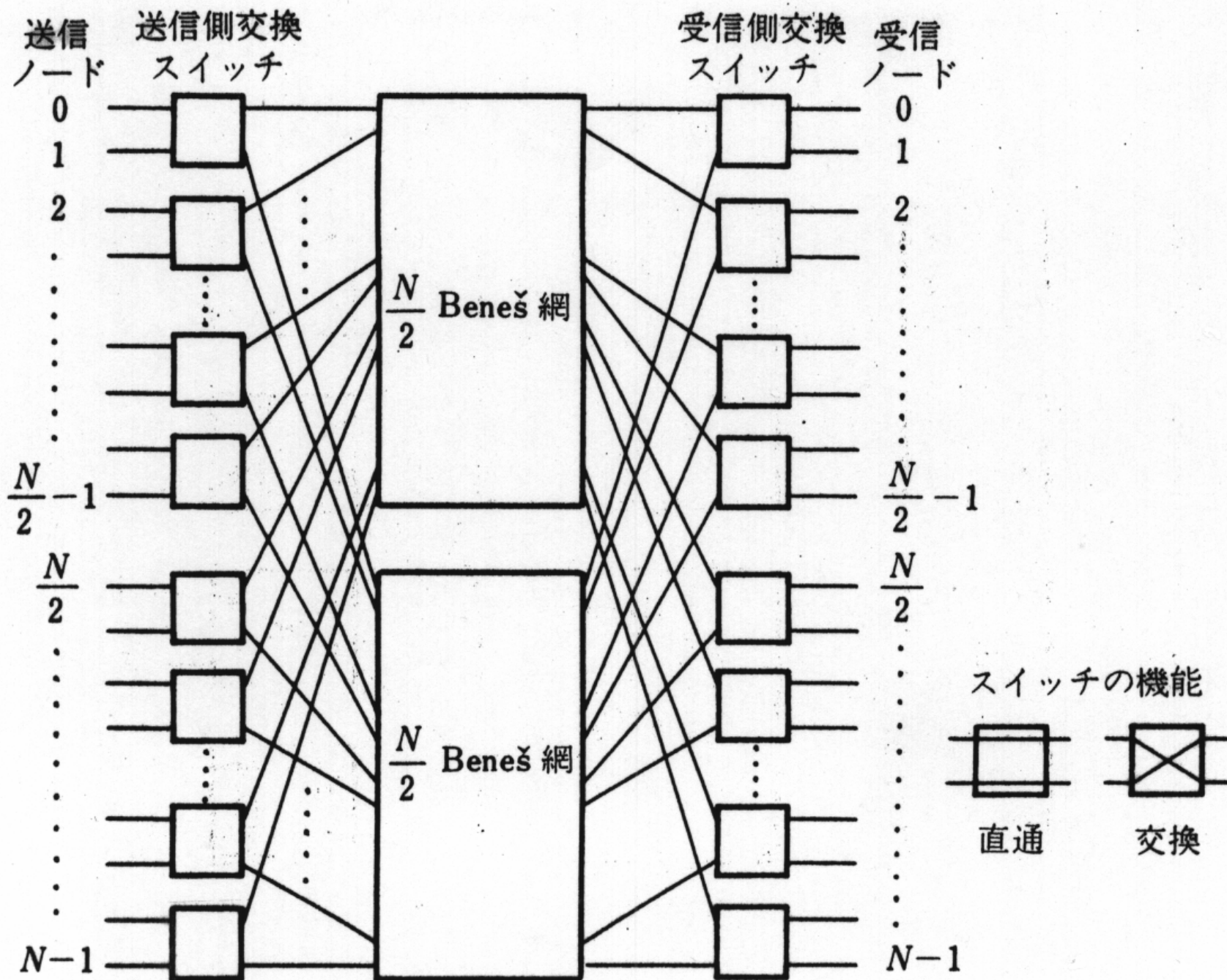
Benes網

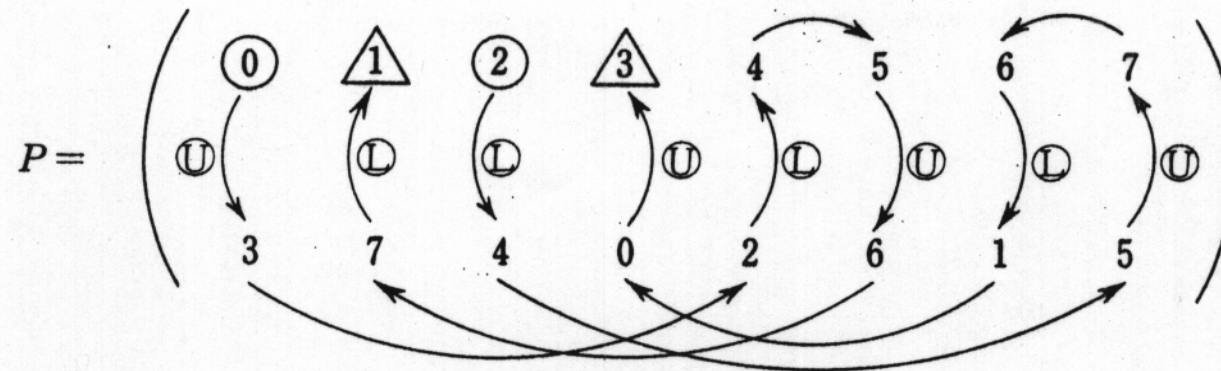
送信ノード	0	1	2	3	4	5	6	7
受信ノード	3	7	4	0	2	6	1	5

経路選択 : ルーピングアルゴリズム

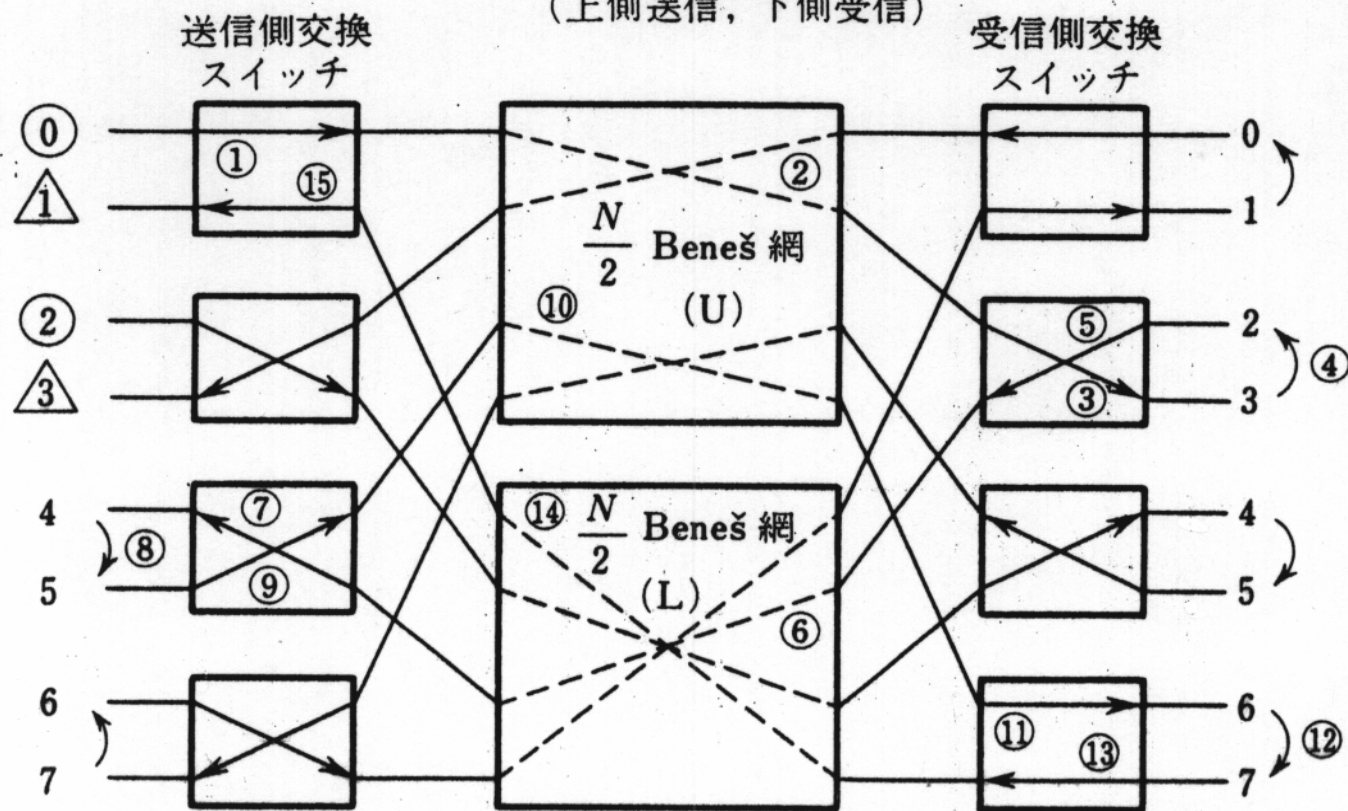


□ : 2×2 スイッチ





(a) 通信パターン  
(上側送信, 下側受信)

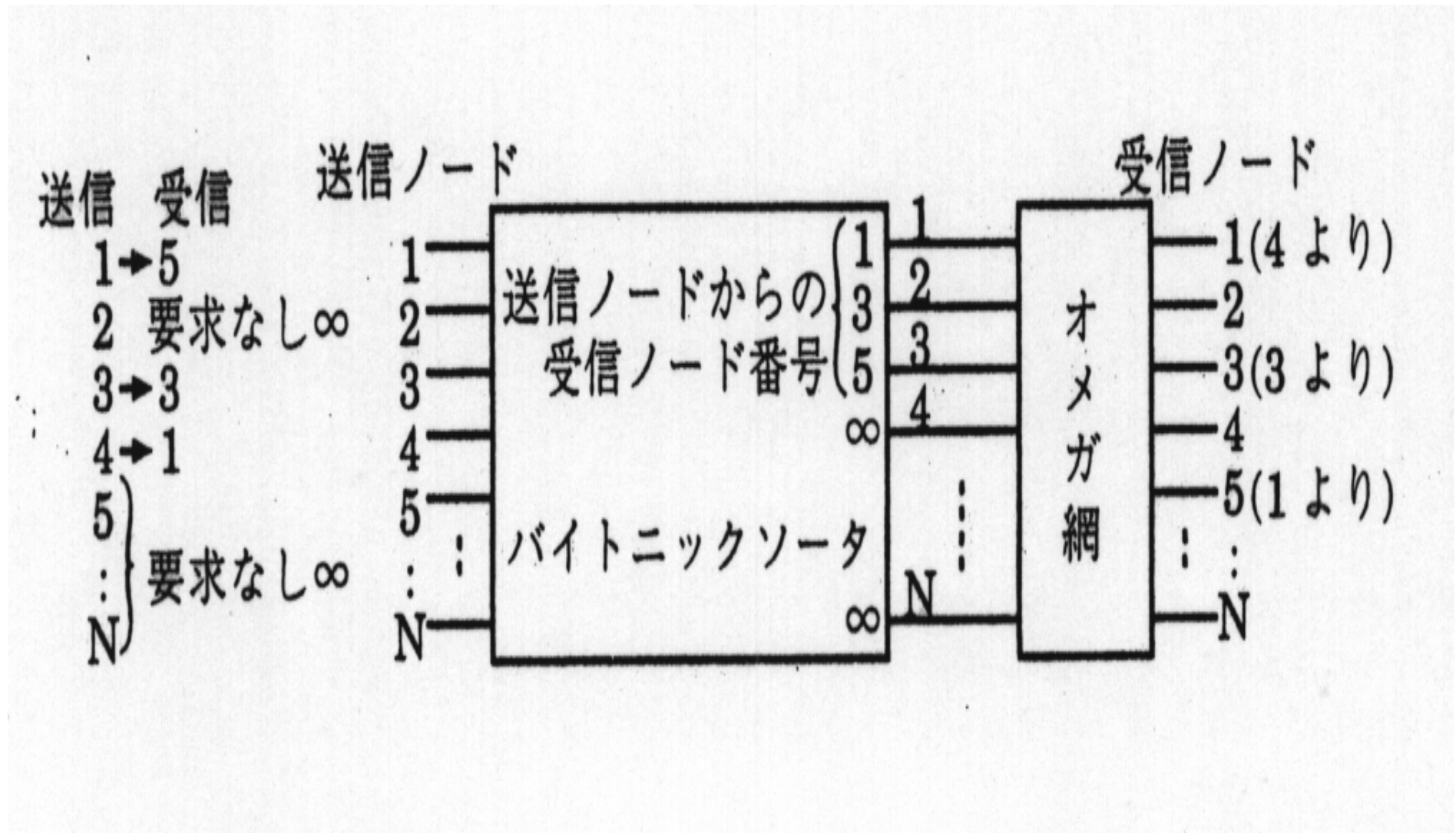


(b) スイッチ設定例



バイトニックソータを用いた結合網

バイトニックソータ（図参照）+ オメガ網



## 非閉塞の証明

$$(a_n, a_{n-1}, \dots, a_2, a_1) > (a_n^*, a_{n-1}^*, \dots, a_2^*, a_1^*)$$

$$(b_n, b_{n-1}, \dots, b_2, b_1) > (b_n^*, b_{n-1}^*, \dots, b_2^*, b_1^*) \text{ なら}$$

$$(a_i, a_{i-1}, \dots, a_2, a_1, b_n, b_{n-1}, \dots, b_{i+1})$$

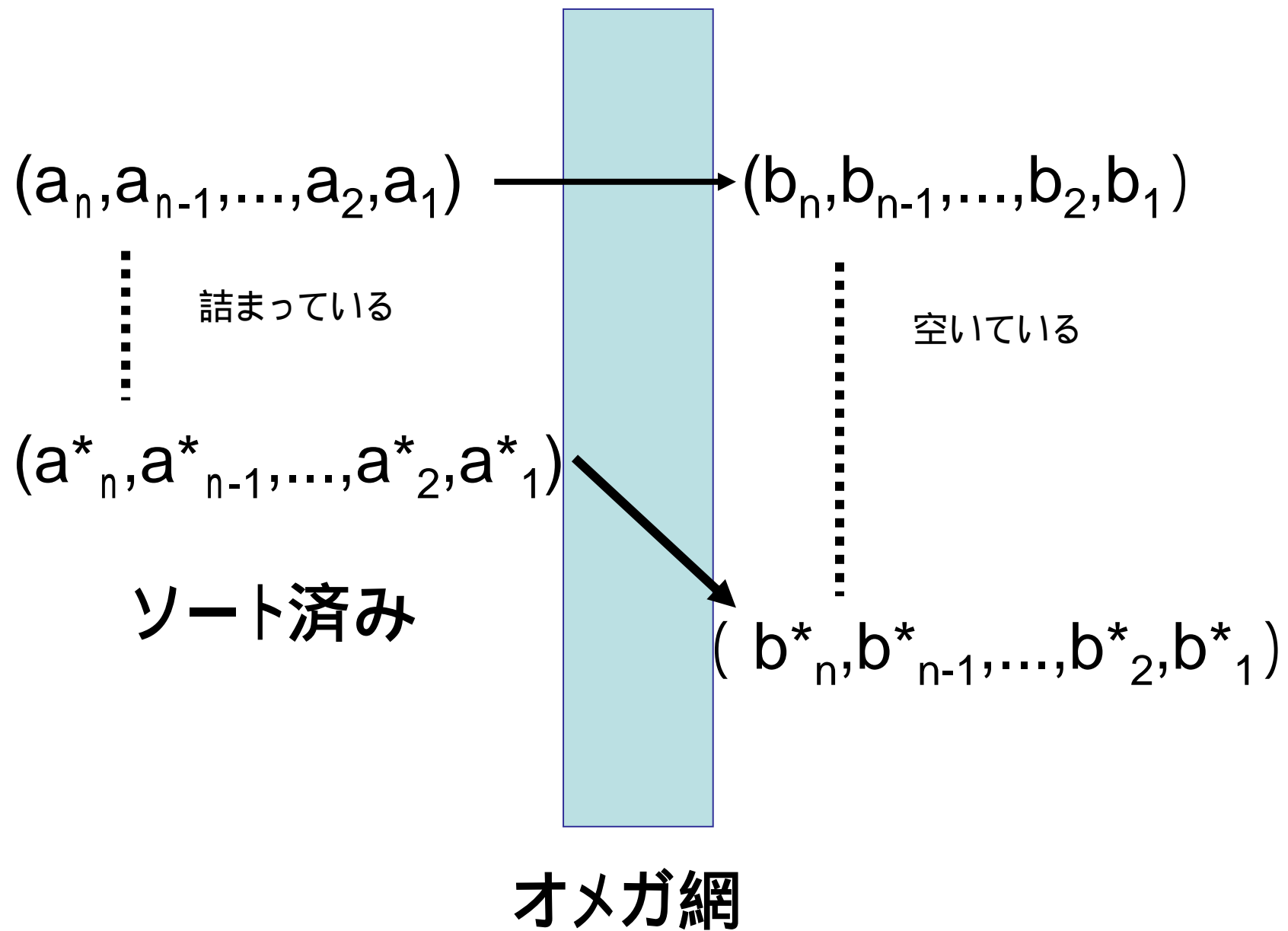
$$(a_i^*, a_{i-1}^*, \dots, a_2^*, a_1^*, b_n^*, b_{n-1}^*, \dots, b_{i+1}^*) \text{ の証明}$$

$$(a_i, a_{i-1}, \dots, a_2, a_1) = (a_i^*, a_{i-1}^*, \dots, a_2^*, a_1^*) \text{ のとき}$$

$$(b_n, b_{n-1}, \dots, b_2, b_1) - (b_n^*, b_{n-1}^*, \dots, b_2^*, b_1^*)$$

$$(a_n, a_{n-1}, \dots, a_2, a_1) - (a_n^*, a_{n-1}^*, \dots, a_2^*, a_1^*) \quad 2^i \text{ なので}$$

$$(b_n, b_{n-1}, \dots, b_{i+1}) > (b_n^*, b_{n-1}^*, \dots, b_{i+1}^*)$$





### 3 ステージ Clos 網

$2n-1 > m$     $n$  : 再構成型非閉塞網

#### Hall の補助定理

$A$  : 集合

$A_1, A_2, A_3, \dots, A_r$  : 部分集合

「 $a_i \in A_i, a_j \in A_j$  ( $i, j=1, 2, \dots, r$ ) で  $a_i \neq a_j$  なる

$A$  の要素が  $A_i, A_j$  に存在する」

「任意の  $k$  個の和集合  $\bigcup_k A_i$  に少なくとも  $k$  個の異なる要素が存在する」

Slepian Duguidの定理

入力ボックス $I_i$ から繋がっている

出力ボックスの集合 $K_i$

$K_i$ は出力ボックス0の部分集合

任意の  $k$  個の部分集合の和  $K_i$  の異なる

要素の数  $t$

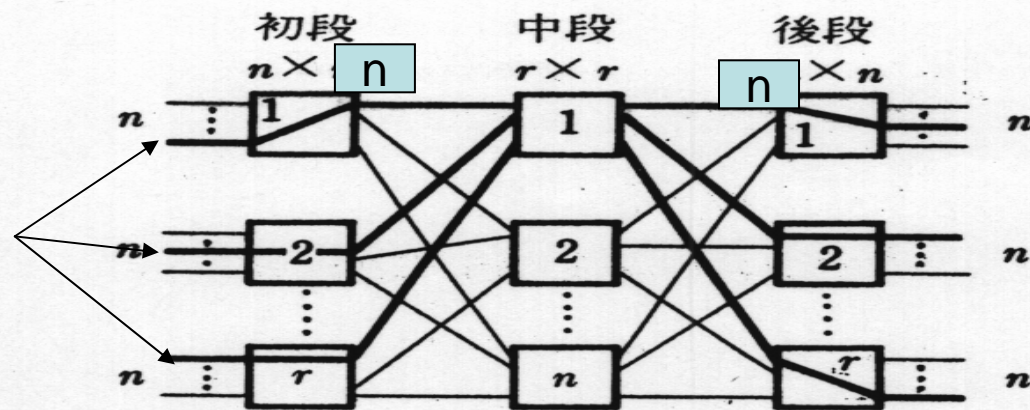
$n$  個の入力ボックスへの入力線数  $k_n$

$K_i$  の出力ボックスからの出力線数  $t_n$

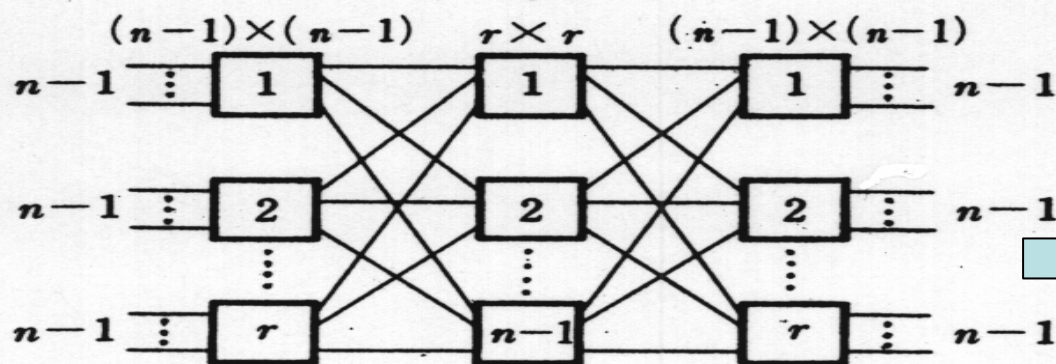
より  $k \leq t$  . Hallの補助定理より

各 $K_i$ には異なる要素が存在する

R本の  
リンク

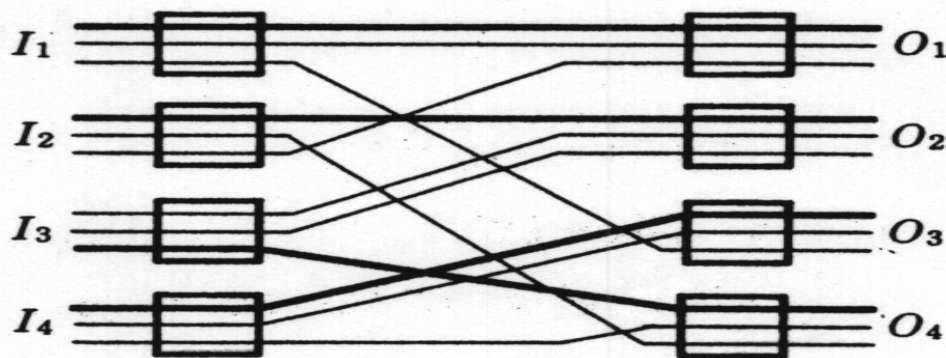
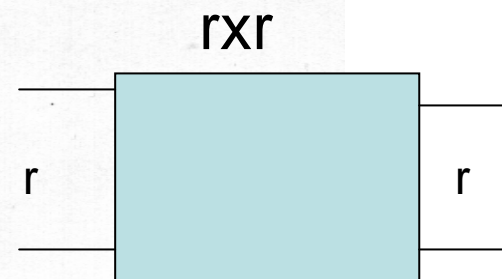
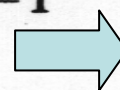


(a)  $\nu(n, n, r)$



(b)  $\nu(n-1, n-1, r)$

与えられた通信パターンに対して  
初段、後段のスイッチボックスを  
1つずつ使用する1本のリンクが  
存在



(c) 証明のための例

$K_1 = (\textcircled{1}, 3)$   
 $K_2 = (1, \textcircled{2}, 4)$   
 $K_3 = (2, \textcircled{4})$   
 $K_4 = (\textcircled{3}, 4)$

### ( 3 ) 閉塞網

可能な通信パターン総数：  $N!$  の一部

$$(\sqrt{N})^N = 2^{(N \log N)/2}$$

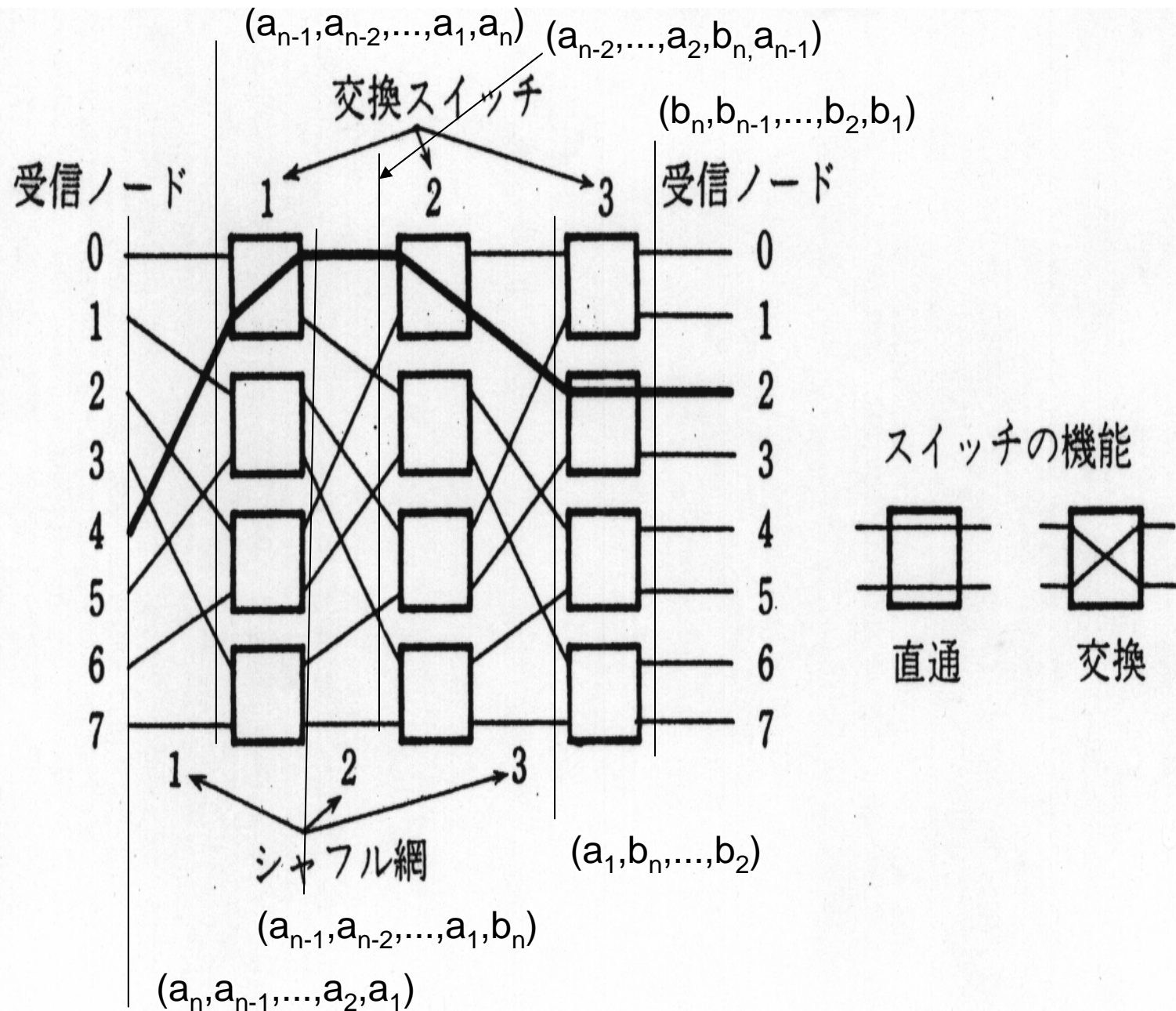
### オメガ網

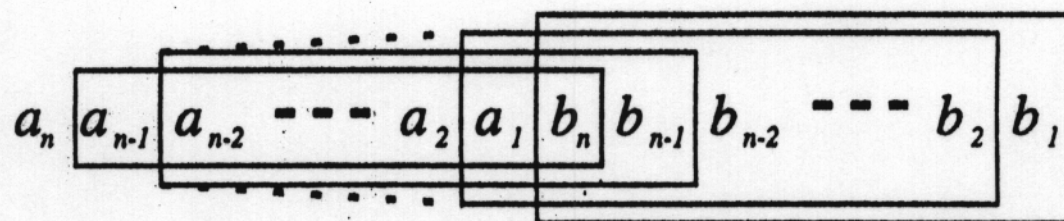
- ・ オメガ網の構成

$$E = (e_1, i)$$

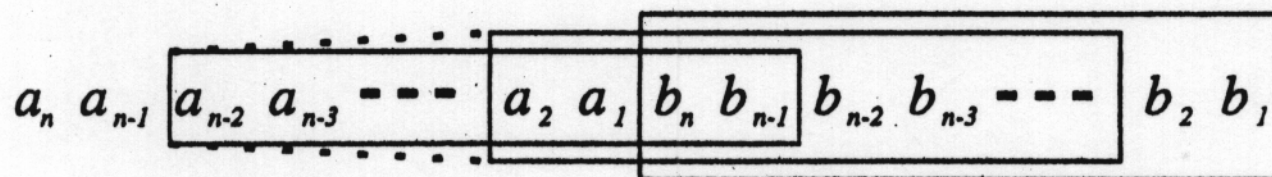
$$\Omega = (\sigma E)^n$$

- ・ オメガ網のスイッチ設定





(a) ウィンドウ (1ビット)



(b) ウィンドウ (2ビット)

送信ノード  $(a_n, a_{n-1}, \dots, a_2, a_1)$

受信ノード  $(b_n, b_{n-1}, \dots, b_2, b_1)$

最初のシャフル：

送信ノード  $(a_n, a_{n-1}, \dots, a_2, a_1)$  は  
 $(a_{n-1}, a_{n-2}, \dots, a_2, a_n)$  に対応

最初の交換スイッチ：

$a_n \oplus b_n = 0$  であれば、直通側

$a_n \oplus b_n = 1$  であれば、交換側

最初の交換スイッチの出力

送信ノード  $(a_n, a_{n-1}, \dots, a_2, a_1)$  は

$(a_{n-1}, a_{n-2}, \dots, a_2, b_n)$

## 第 2 段のシャフル

$(a_{n-2}, a_{n-3}, \dots, a_2, b_n, a_{n-1})$

## 第 2 段の交換スイッチ：

$a_{n-1} \oplus b_{n-1} = 0$  のとき、直通

$a_{n-1} \oplus b_{n-1} = 1$  のとき、交換

## 第 2 段交換スイッチの出力

送信ノード  $(a_n, a_{n-1}, \dots, a_2, a_1)$

$(a_{n-2}, a_{n-3}, \dots, a_2, b_n, b_{n-1})$  に対応

交換スイッチ



$b_i=0$ なら、上方リンク選択

$b_i=1$ なら、下方リンク選択

ウィンドウ表現

BxBスイッチの利用

$$B=2^b$$

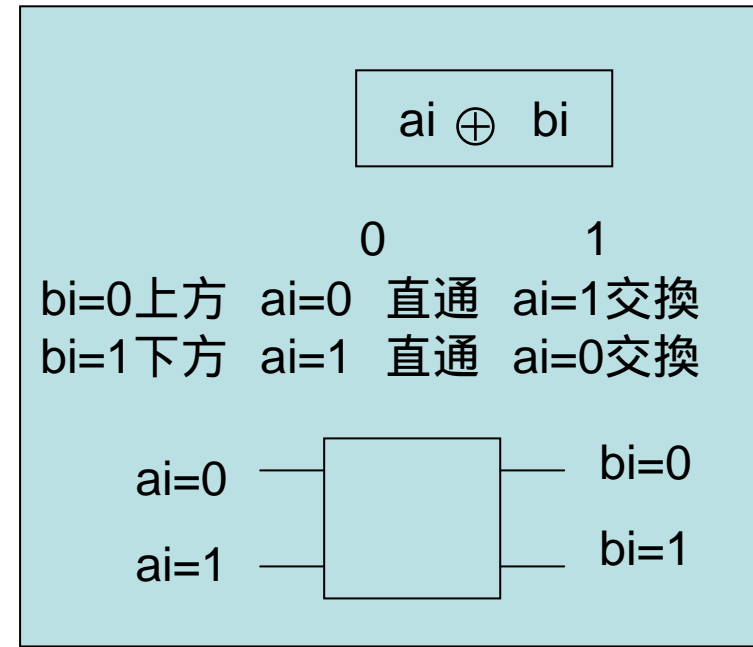
$$N=B^k$$

BxBのスイッチ k 段のオメガ網

2 ビットシャフル、

$(a_n, a_{n-1}, \dots, a_2, a_1)$

$(a_{n-2}, a_{n-3}, \dots, a_2, a_1, a_n, a_{n-1})$

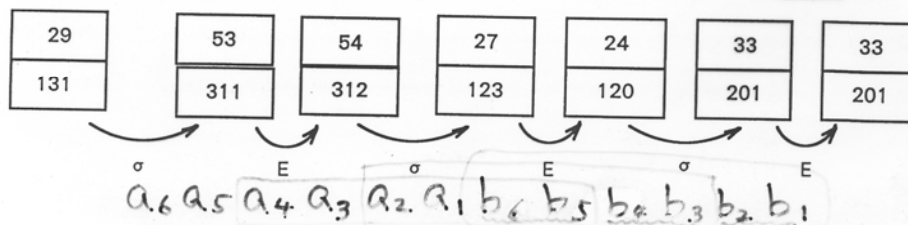


# 64入力-64出力オメガ網

011101

100001

011101 100001



## 間接 2 進 n - キューブ

$$C = E\beta_2 E\beta_3, \dots, E\beta_n E\sigma^{-1}$$

## バンヤン (Banyan) 網

$$Y = E\beta_2 E\beta_3, \dots, E\beta_n E$$

## ベースライン (baseline) 網

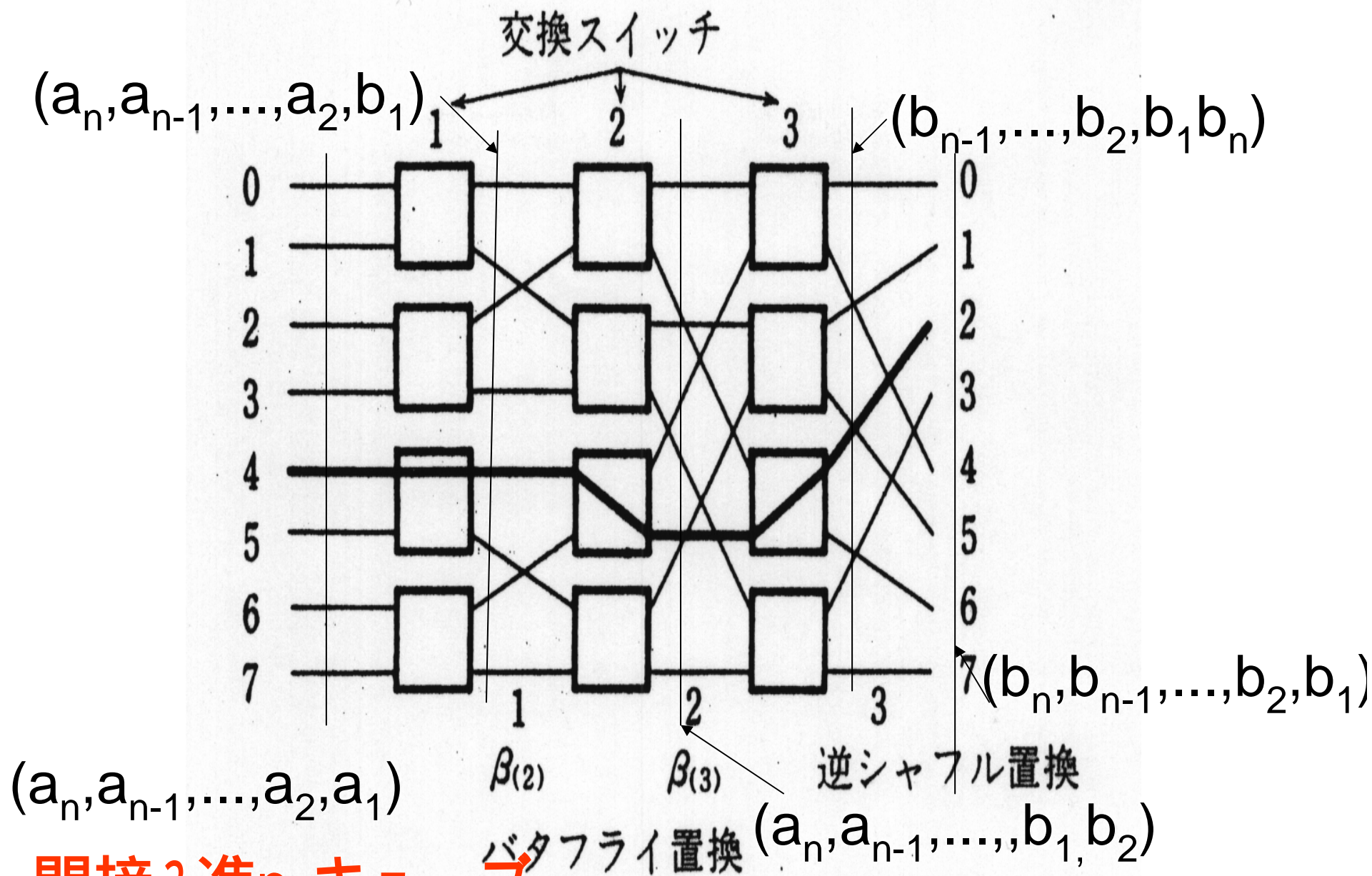
$$B = E\sigma_n^{-1} E\sigma_{n-1}^{-1}, \dots, E\sigma_2^{-1} E$$

## 閉塞網間の関係

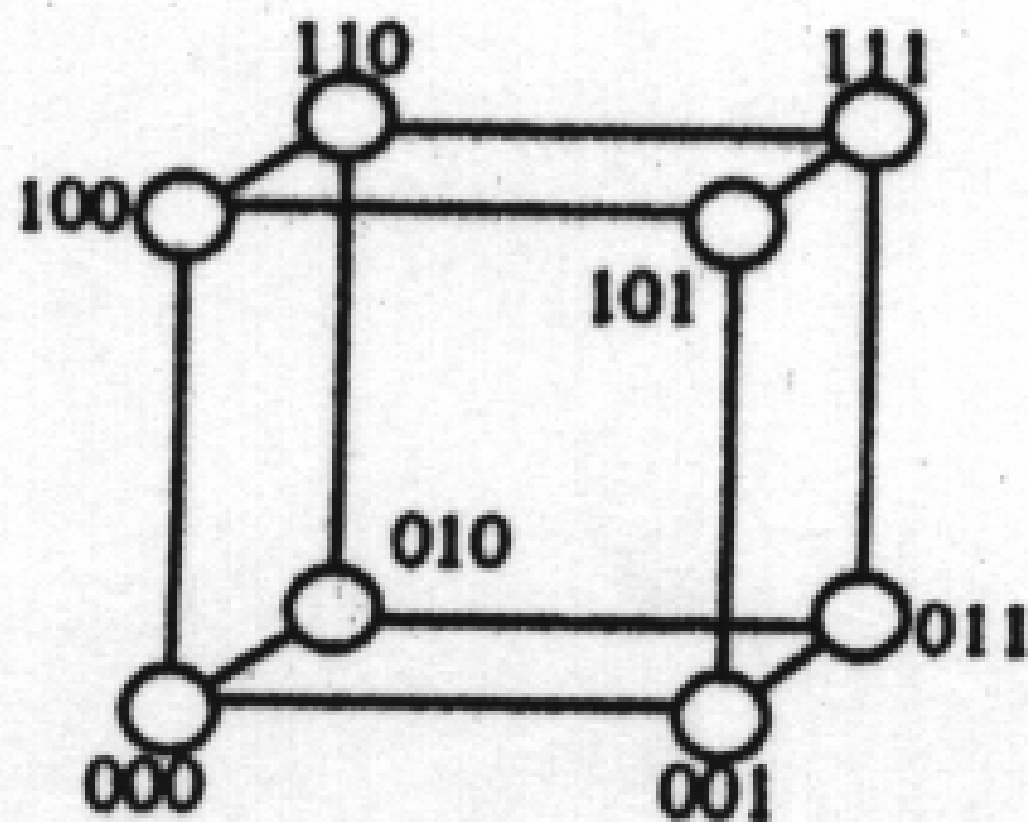
オメガ網 (  $\Omega$  )、間接 2 進 n キューブ網 ( C ) ,

バンヤン網 ( Y )

$$\Omega^{-1} = C = Y\sigma$$



間接2進 $n$ -キューブ



(a) 2進 3-キューブ

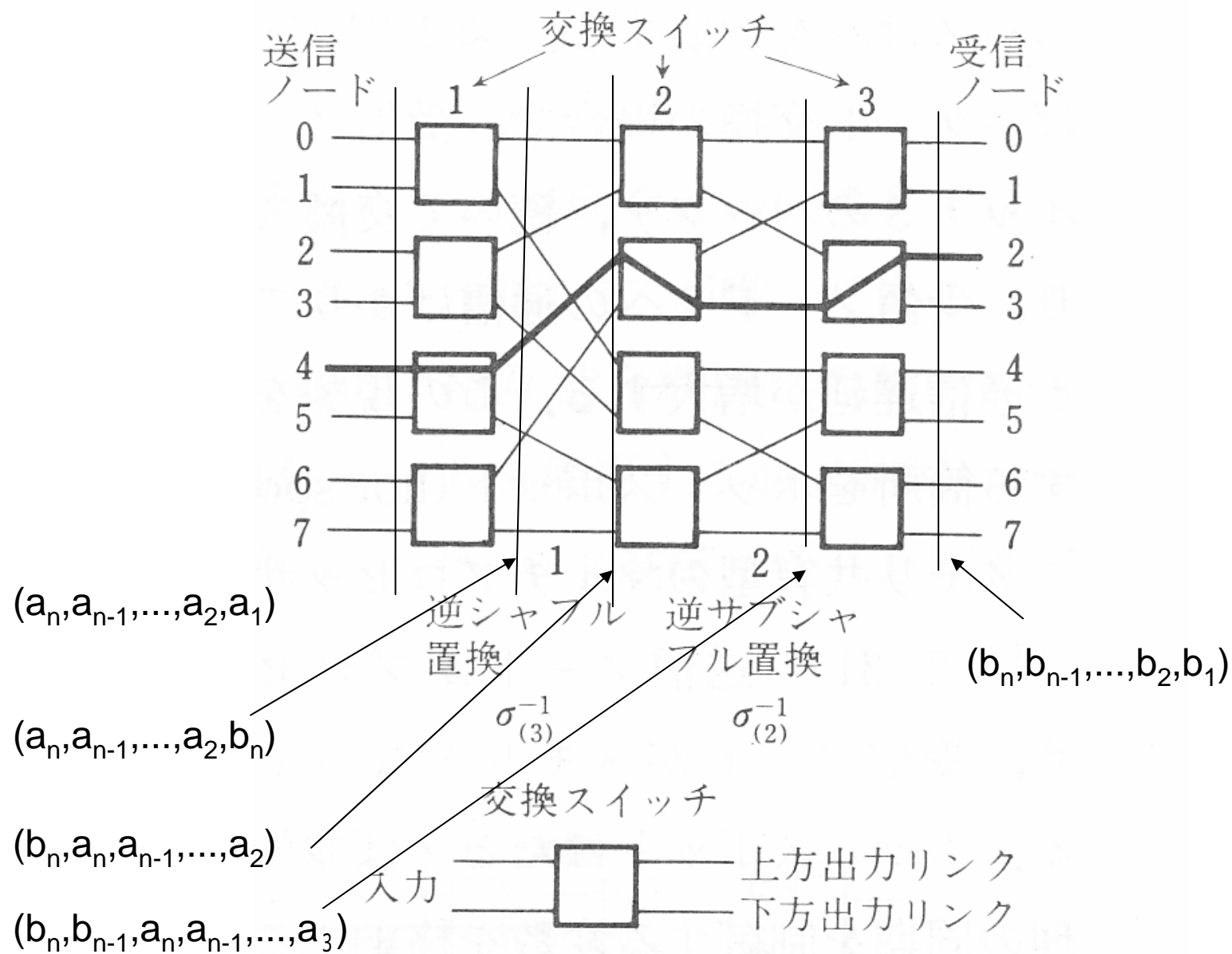
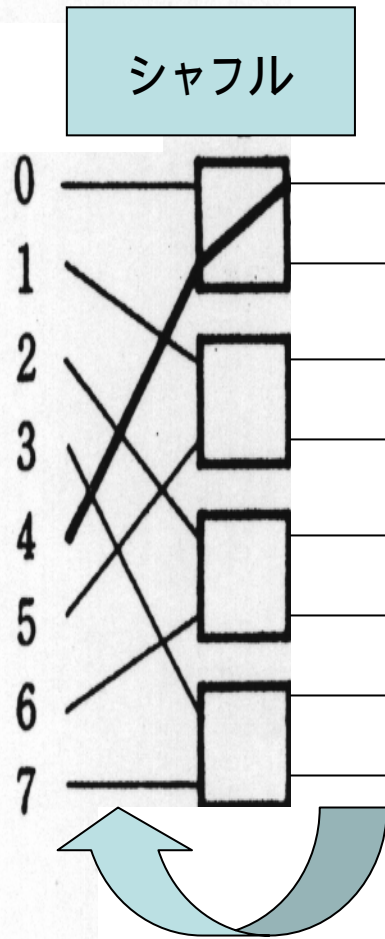


図 2.30 ベースライン網

# 単一段結合網



## 単一段シャフルネットワーク

## 2.5 多段結合網の通信 制御と耐故障設計

### 2.5.1 木飽和と結合操作

例 ベクトル処理

プロセッサ  $P_0, \dots, P_7$  :

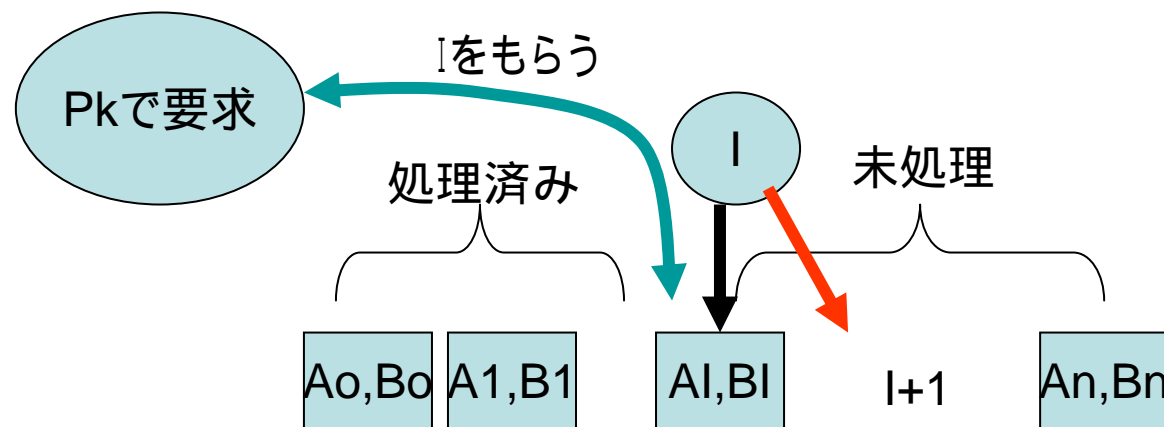
$A(I) + B(I)$

$I$  を獲得し,  $+1$

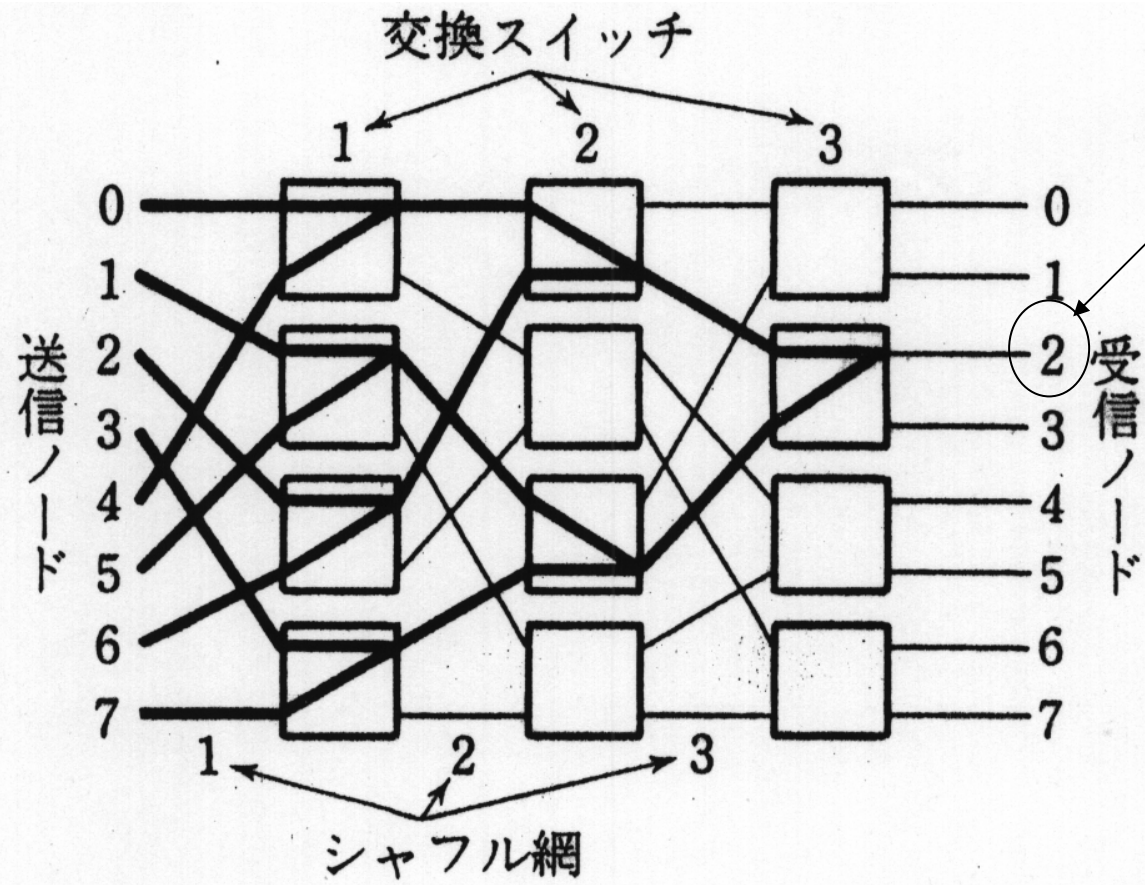
共有変数  $I$  : ホットスポット

Fetch and Add 命令

あるスイッチで2つの要求







制御変数格納

F&A 1      F&A 1 - > F & A2  
    ↖ 10      ↖ 11      I=10      12  
                コンバイニング

## 2.5.2 耐故障網

送信側に付加ステージを追加する方式

各ステージ内にループ構造を導入する方式

大きなクロスバスイッチを用いる方式

クロスバスイッチの大きさ： $B \times B$  ( $B=2^b$ )、

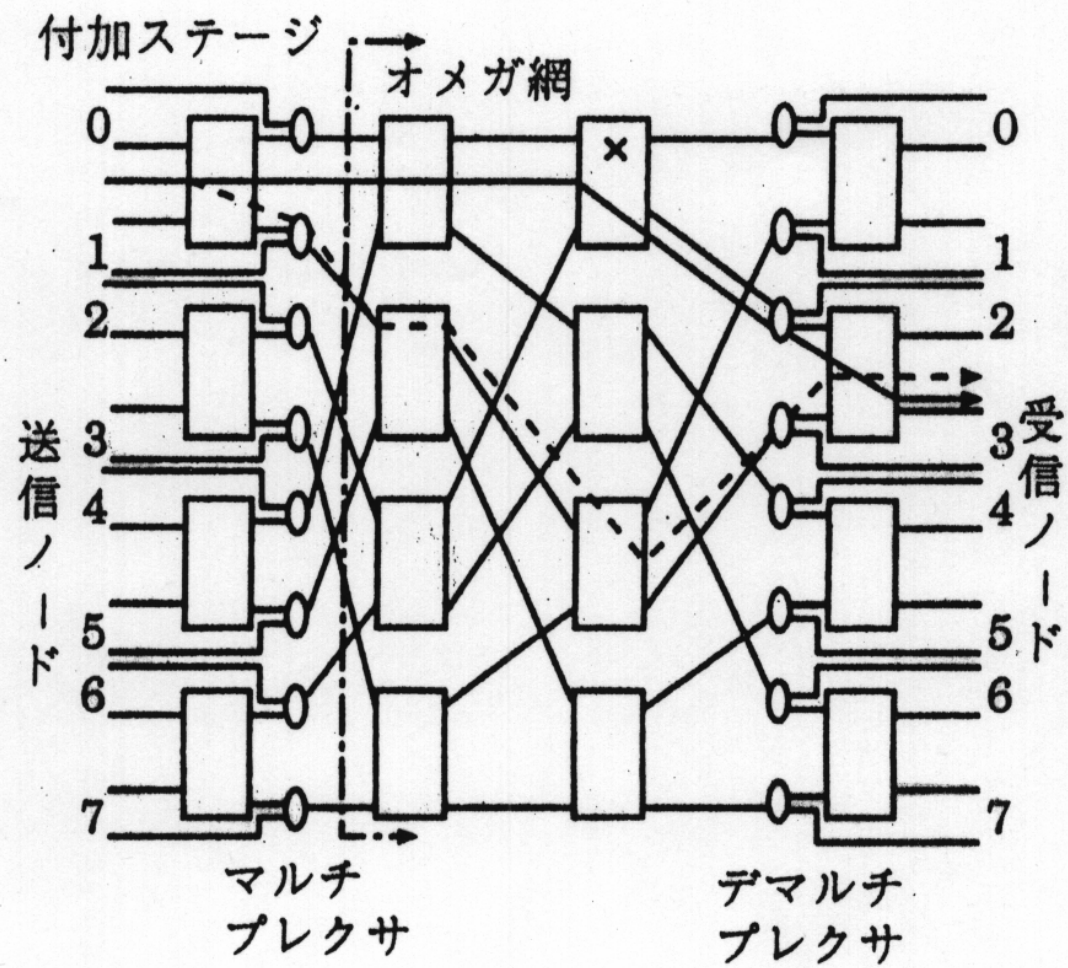
ノード数  $N (=2^n)$

R - パス可能な  $\lceil \log_B N \rceil$  段の

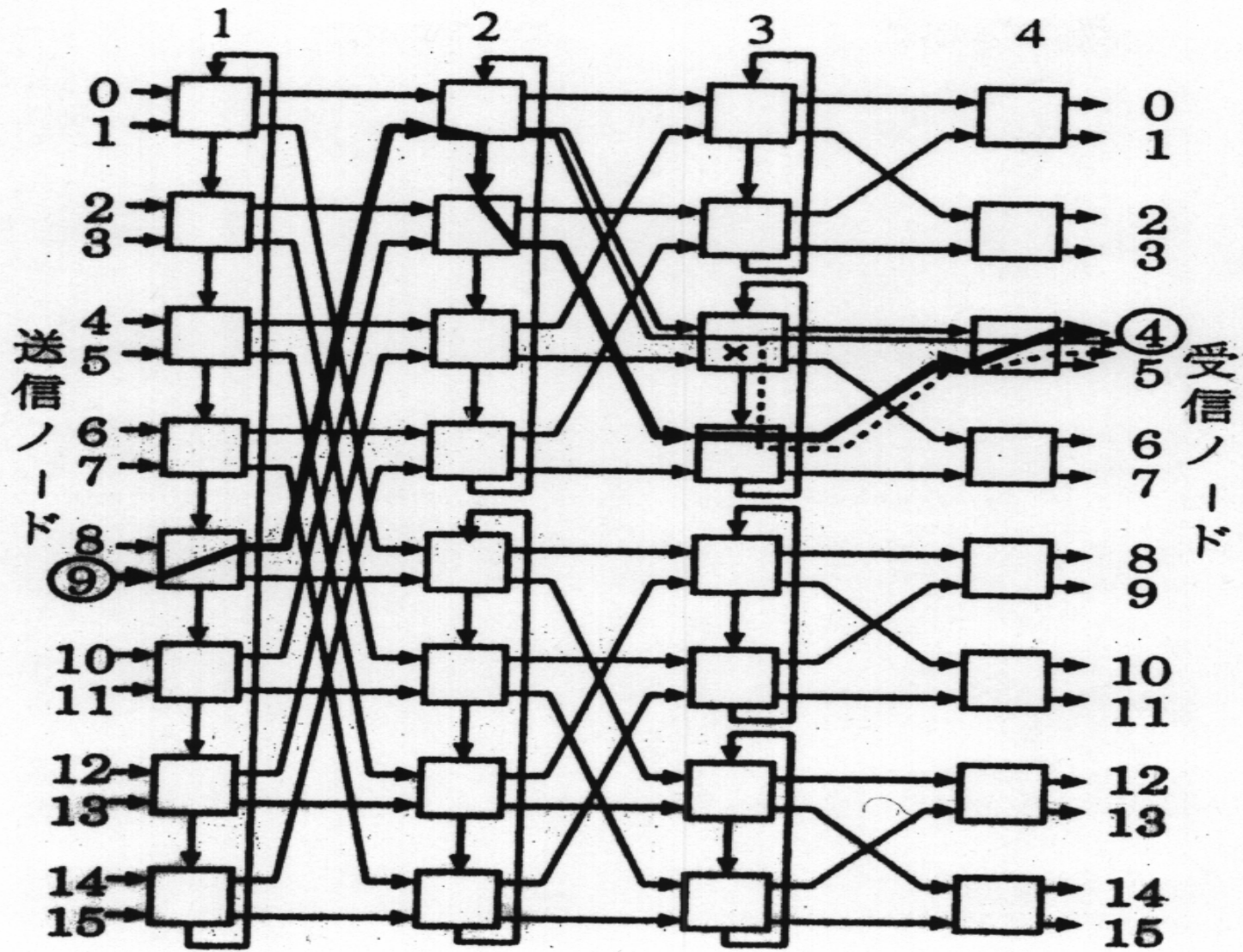
マルチパスオメガ網

$$R = B^{\lceil n/b \rceil - n/b}$$

$N=2^n=B^k$  のとき、 $R=1$



# ステージ



## ( 1 - パスオメガ網 )

### k の設定

$$2^{b(k-1)} < N < 2^{bk}$$

$$k = \lceil \log_B N \rceil = \lceil n / b \rceil$$

kb-n個の \*

$$a_n a_{n-1} \cdots a_{n-b} \cdots a_{n-2b} \cdots a_1 \text{****} b_n b_{n-1} \cdots b_2 b_1$$

長さ n ビットのウィンドウ : b ビット毎、 k 個

$$kb - n < b$$

最終ステージ以外のウィンドウ :

kb-n個の \* 印をすべて含む。

$$R = 2^{kb-n} = B^{\lceil n/b \rceil - n/b}$$

N=32、B=4の3ステージのオメガ網

ノード8から12への通信

ウィンドウ 0 1 0 0 0 \* 0 1 1 0 0

\* = 0 の時

各ステージでの受信ノード番号

0 0 0 0 0、0 0 0 1 1、0 1 1 0 0

経路 $e_1$ 、 $e_2$ 、 $e_3$ 、 $e_6$

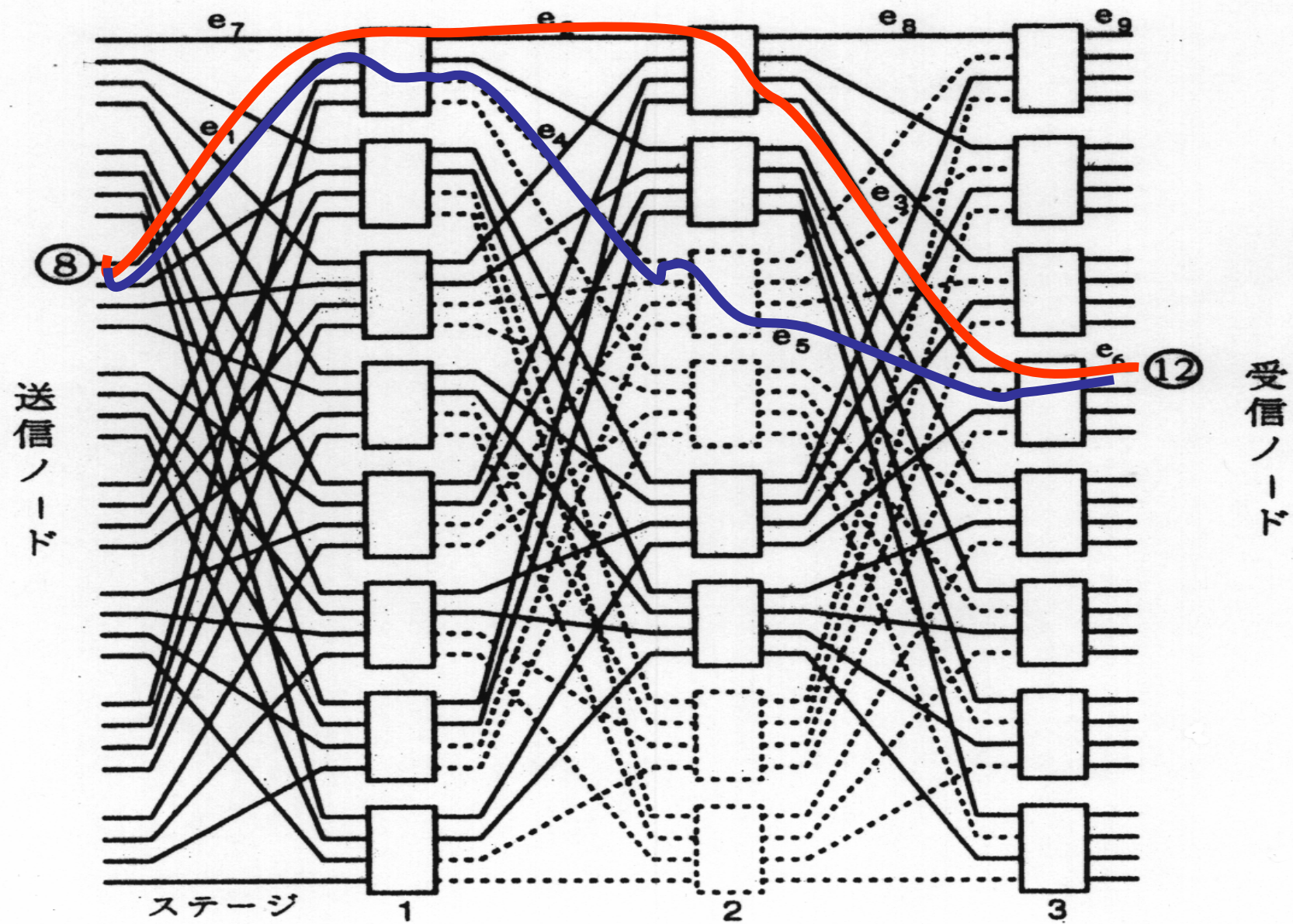
\* = 1 の時

各ステージでの受信ノード番号

0 0 0 1 0、0 1 0 1 1、0 1 1 0 0

経路 $e_1$ 、 $e_4$ 、 $e_5$ 、 $e_6$







## 2.5.3 負荷分散網

### 1.6 可変構造型

#### 相互結合網

機械的可変構造

可変トーラス網

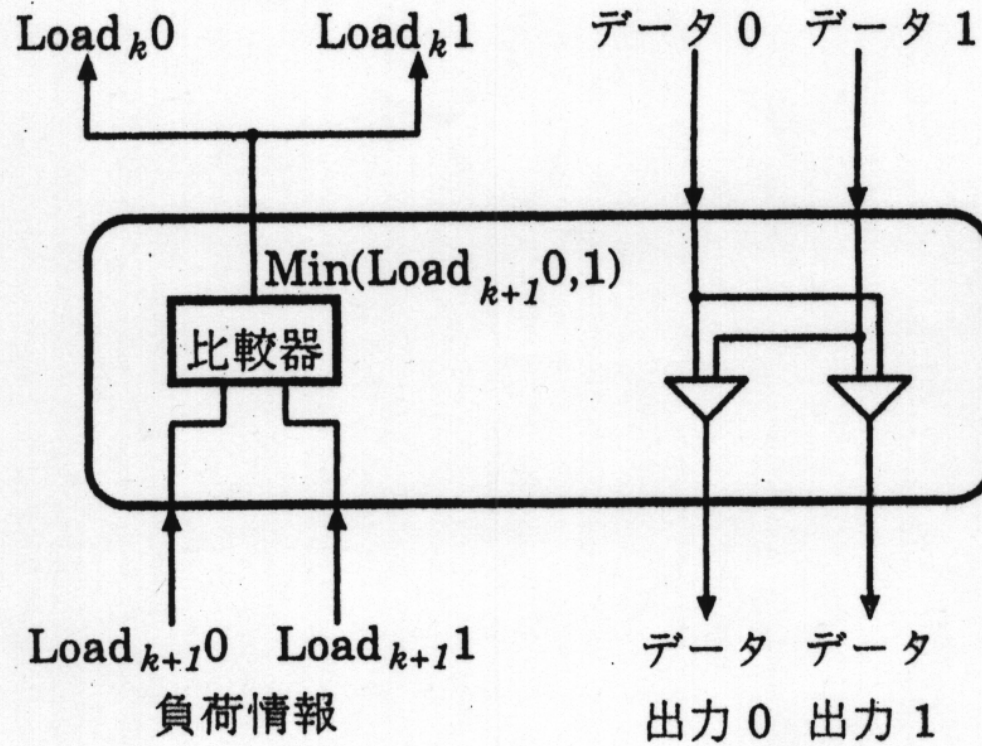
多重クロスバ網

高速網シミュレーション

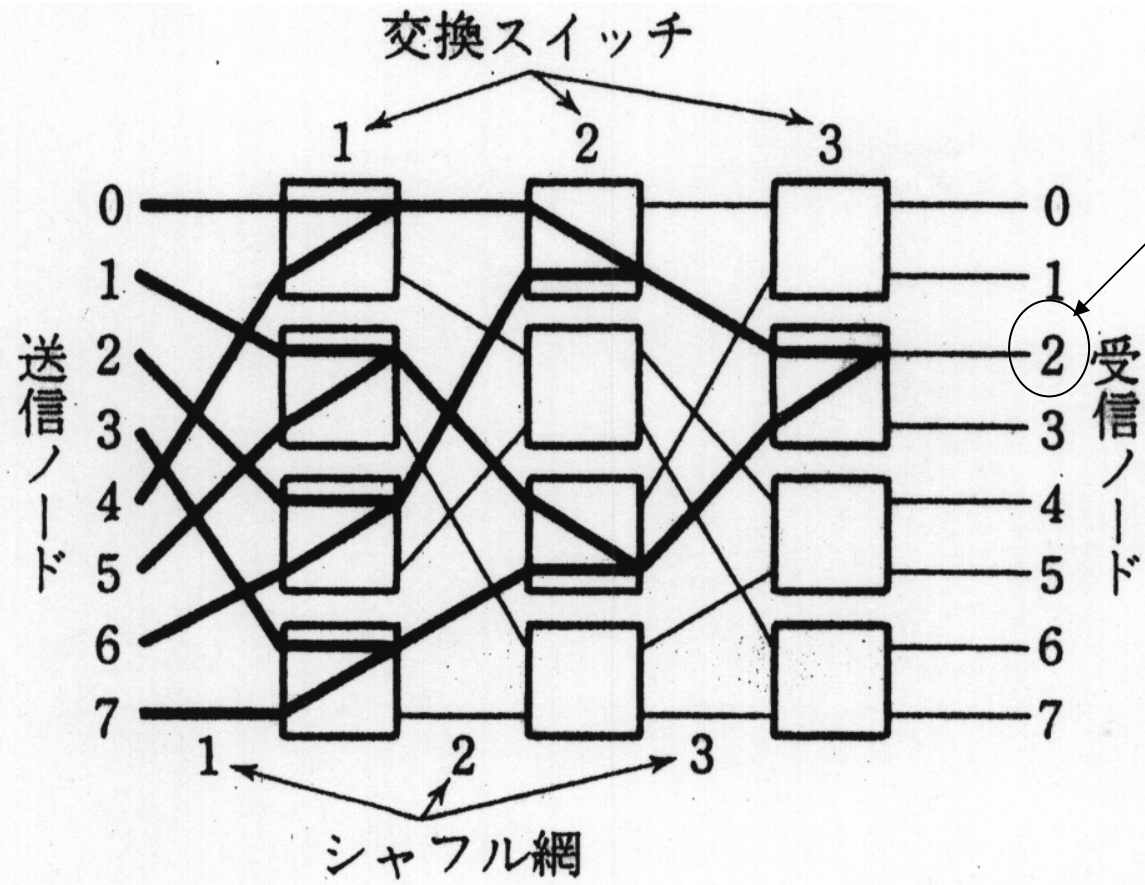
動的リンク形成

送信プロセッサ側

負荷情報

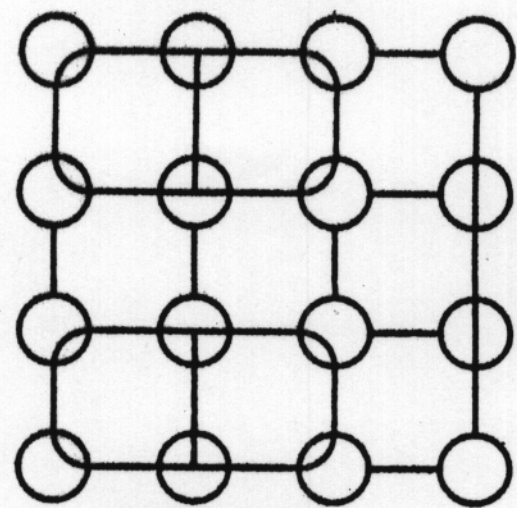


受信プロセッサ側



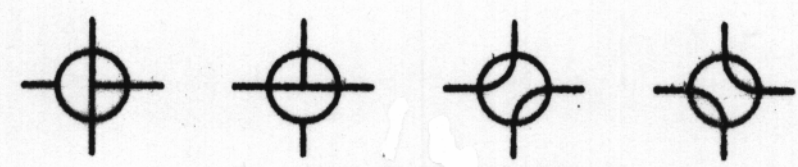
制御変数格納

トーラス1

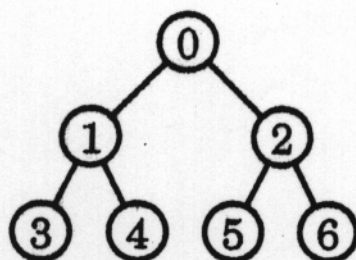


リング

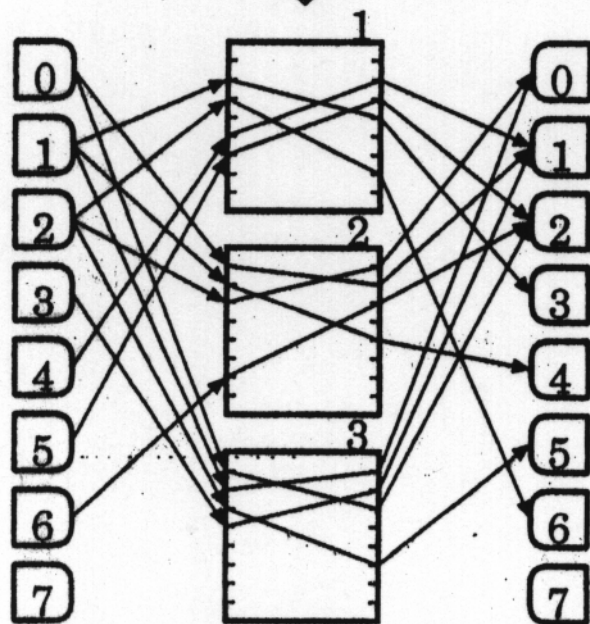
トーラス2



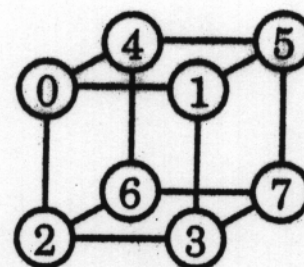
スイッチ内での結合



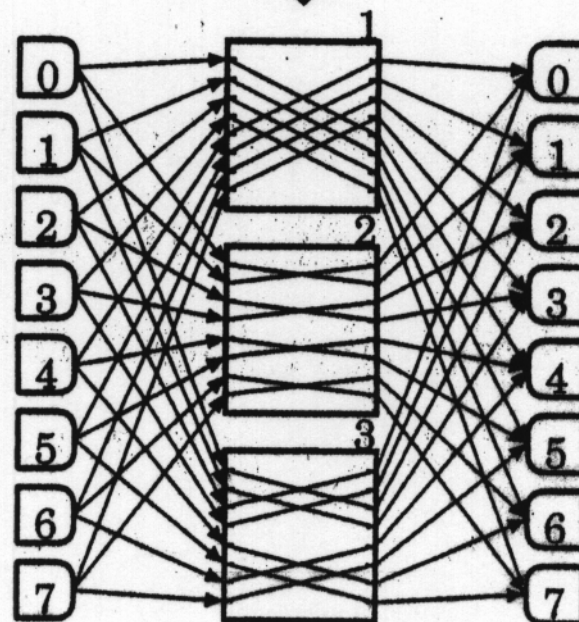
結合パターン



(a) トリー網



結合パターン



(b) パイパキューブ網

1, 2, 3 : 通信パターン