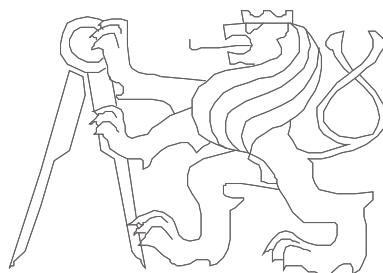


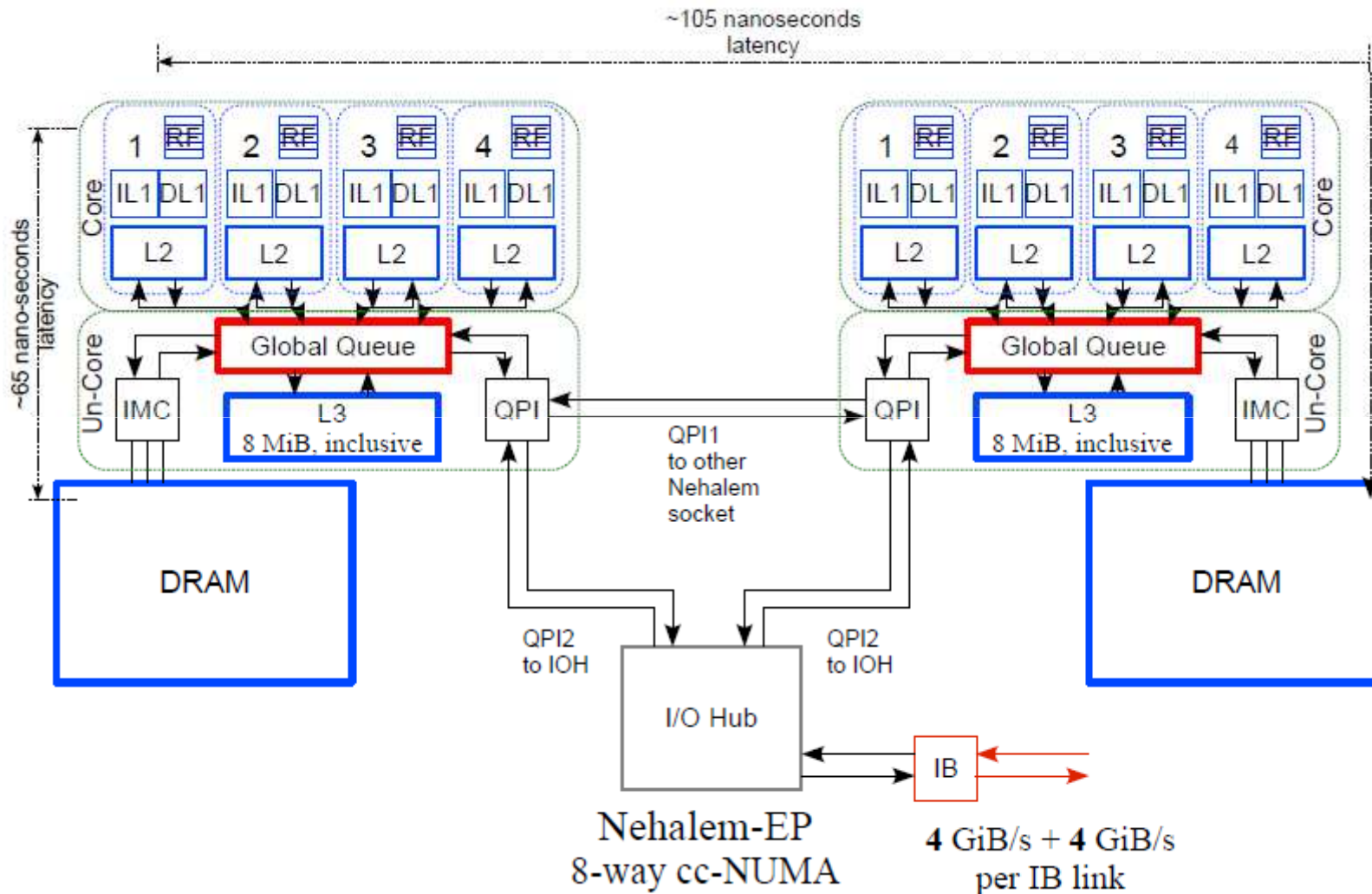
# Pokročilé architektury počítačů

## Propojovací sítě

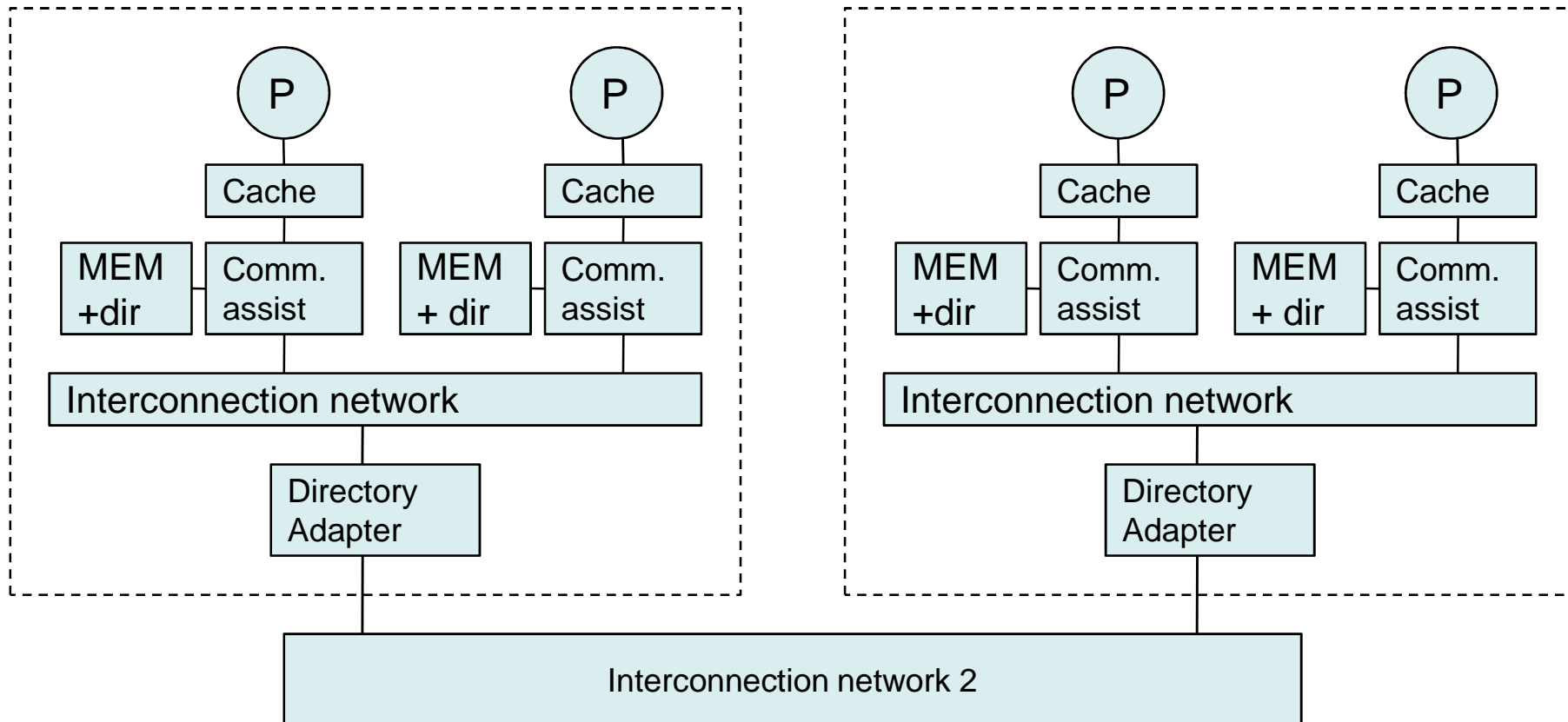


České vysoké učení technické, Fakulta elektrotechnická

# Motivace: Propojení dvou procesorů pomocí QPI



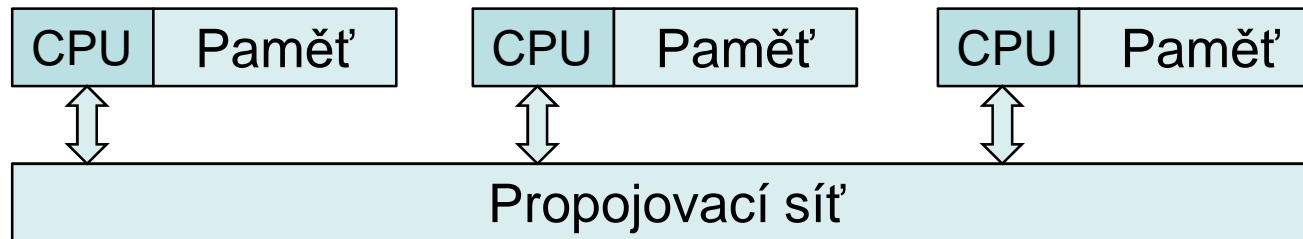
## Two-level cache coherent system



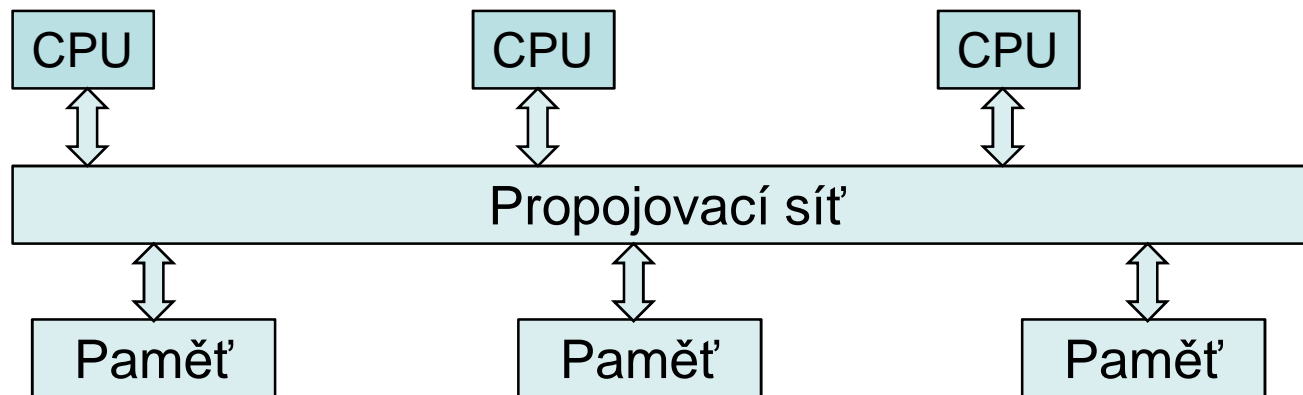
- Directory-directory
- Alternatives: Snooping-Snooping, Snooping-Directory, Directory-Snooping

## Paměťové architektury paralelních procesorů

- **Volně vázaný systém** (Loosely coupled) – paměť je distribuována mezi uzly a každý uzel může přistoupit k paměti v jiném uzlu.  
Pozn.: Pokud by tuto možnost neměl, pak se nejedná o model se sdílenou pamětí.



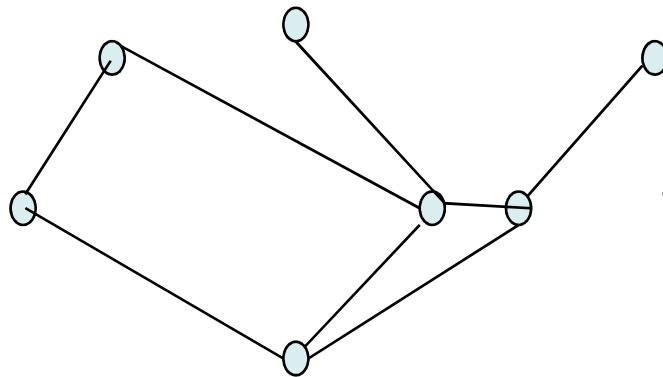
- **Těsně vázaný systém** (Tightly coupled) – paměť je lokalizována centrálně. Procesor může obsahovat vlastní paměť/cache.



## SW pohled na propojovací síť

Síť musí/by měla podporovat:

- One-to-one komunikaci
- One-to-all broadcast a All-to-one redukci
- All-to-all broadcast a All-to-all redukci
- Scatter a Gather



Jak vypadá síť, která optimalizuje tuto komunikaci??

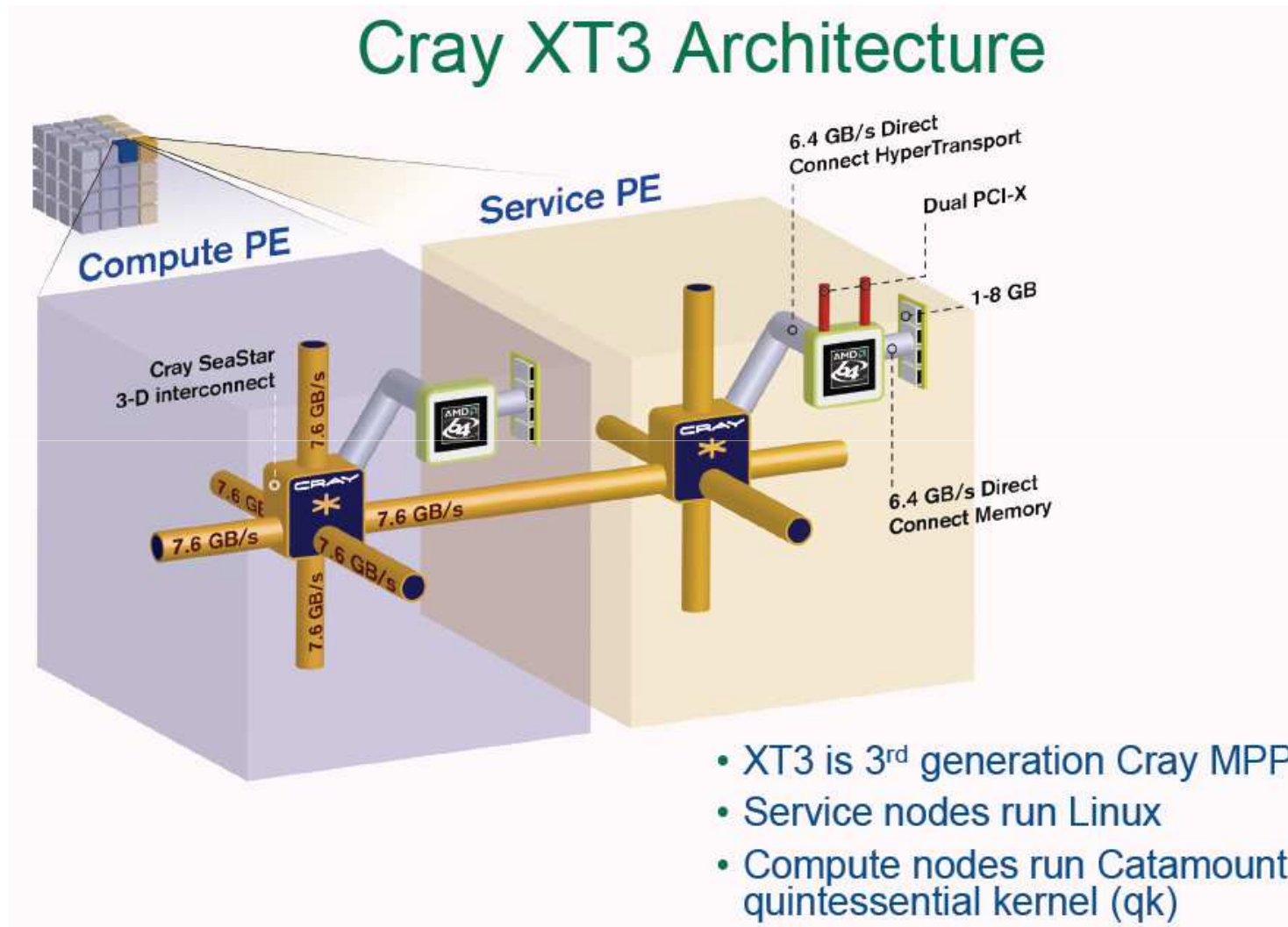
## Motivace - Cray XT3

Cray XT3 (20.5 TFlops) v Oak Ridge National Laboratory:

- 56 cabinets, 5212 výpočtových procesorových elementů (PEs), 82 servisních procesorových elementů
- výpočtový uzel – je tvořen 4 CPU (2GB/PE)
- CPU: 64-bitový 2.4GHz AMD Optreon



# Motivace - Cray XT3



# Topologie

## Topologie propojovacích sítí

### Statické sítě

- lineární pole (linear array)
- kruh (ring)
- Chordálový kruh (Chordal ring)
- binary tree
- fat tree
- 2D, 3D mesh
- 2D, 3D torus
- hypercube
- Cube Connected Cycles (CCC)

### Dynamické sítě

- Bus network
- jednostupňové sítě (křížové přepínače - crossbars)
- vícestupňové sítě (omega, Banyan, Cantor, Clos,...)

## Topologie propojovacích sítí

### Přímé sítě

Každý uzel má svůj přepínač a obráceně

### Nepřímé sítě

Některé přepínače nemají přiřazen uzel, pouze směřují provoz k dalším přepínačům



# Statické sítě

## Parametry

- **velikost sítě  $N$**  (Network size): počet uzlů v síti,
- **stupeň uzlu  $d$**  (Node degree): počet hran vstupujících nebo vystupujících z uzlu,
- **bisekční šířka  $B$**  (Bisection width): minimální počet hran, které musíme přerušit při rozdělení sítě na dvě stejné poloviny,
- **průměr sítě  $D$**  (Network diameter): počet hran maximální nejkratší cesty mezi dvěma libovolnými uzly sítě - ukazuje nejdelší komunikaci,
- **cena  $C$**  (Cost): počet komunikačních linek (hran)



**Linear array**

# Statické sítě



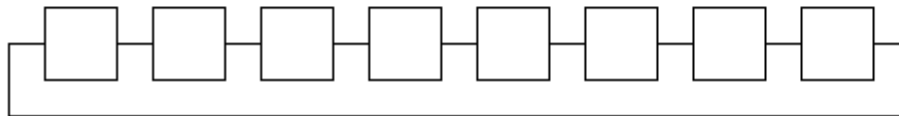
**Linear array**

**d:** 1 (pre koncové uzly), 2 (pre ostatné)

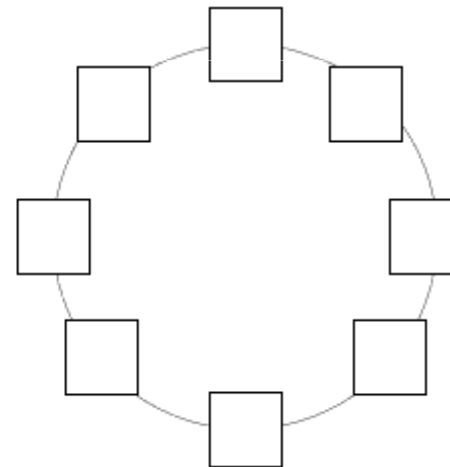
**D:** N-1

**B:** 1

**C:** N-1



**Ring**



**d:** 2

**D:** N/2

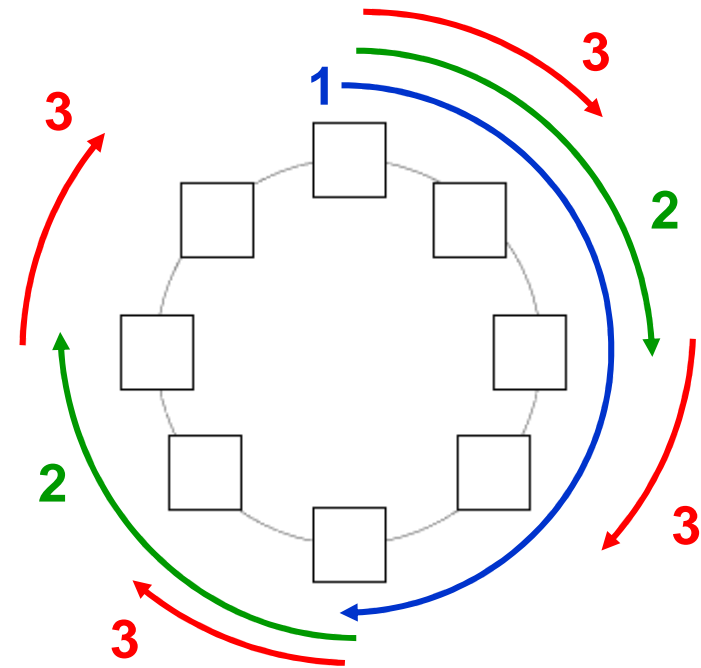
**B:** 2

**C:** N

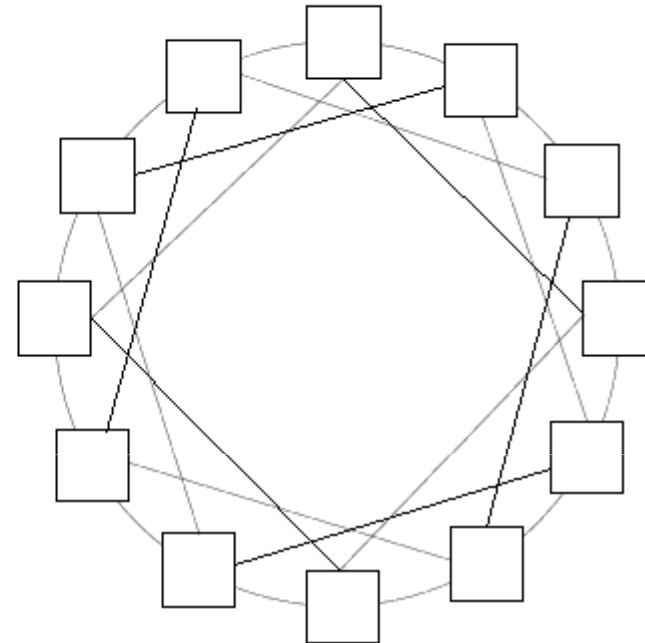
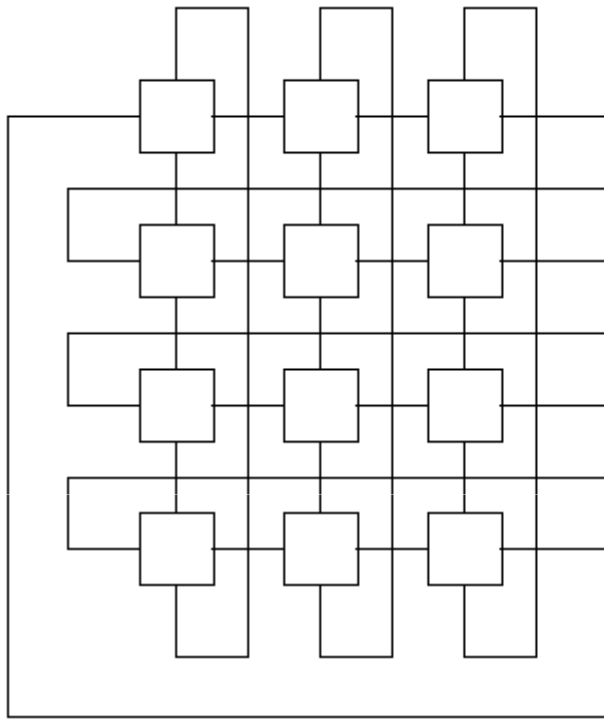
## Ring – příklad komunikace

### One-to-all broadcast a All-to-one redukce

- Nejjednodušší je poslat N-1 zpráv ze zdroje do N-1 cílů (procesorů, uzlů). Není moc efektivní...
- Použij algoritmus rekurzivní zdvojení (recursive doubling) – zdroj pošle zprávu do vybraného uzlu. Tím se problém rozdělí na dva.
- Nebo lineárně sousedovi
- Redukce – obráceně. Data musí však zpracovat...

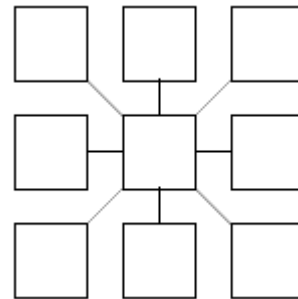


## Statické sítě



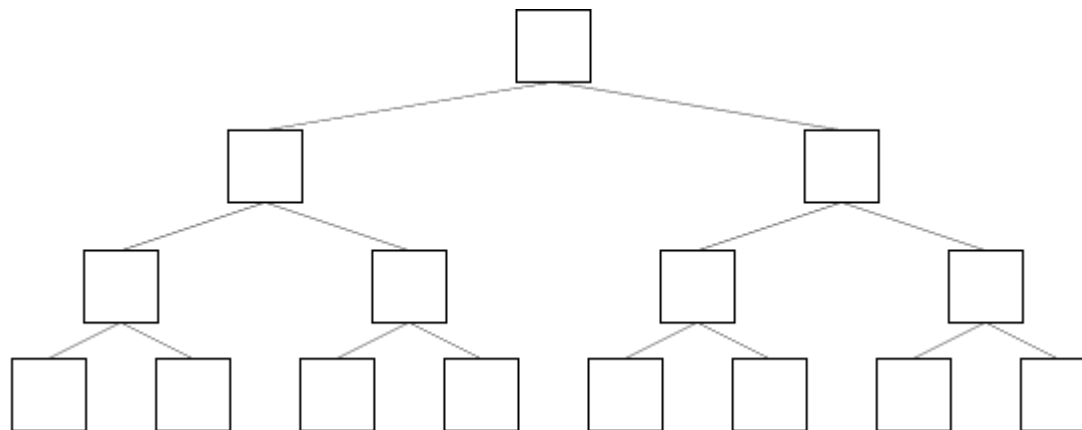
**Chordal ring**

# Statické sítě



**d:** 1, N-1  
**D:** 2  
**B:** 1  
**C:** N-1

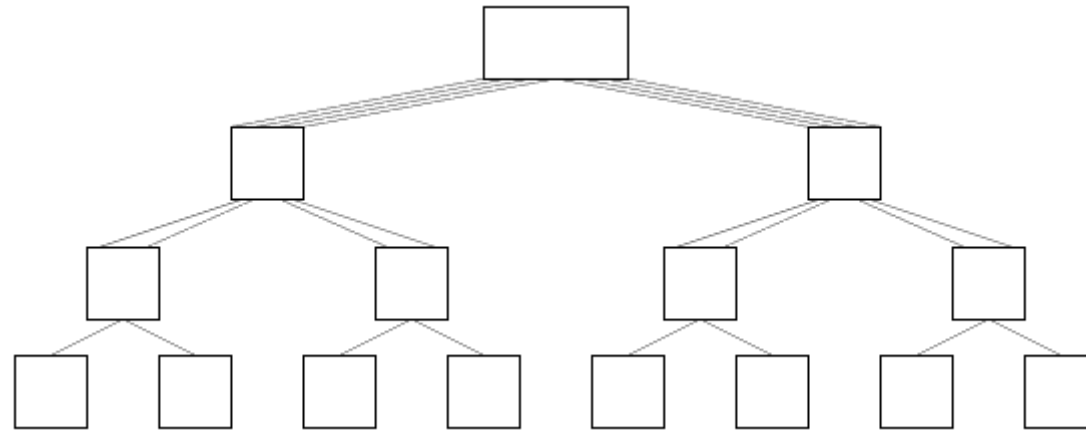
**Star**



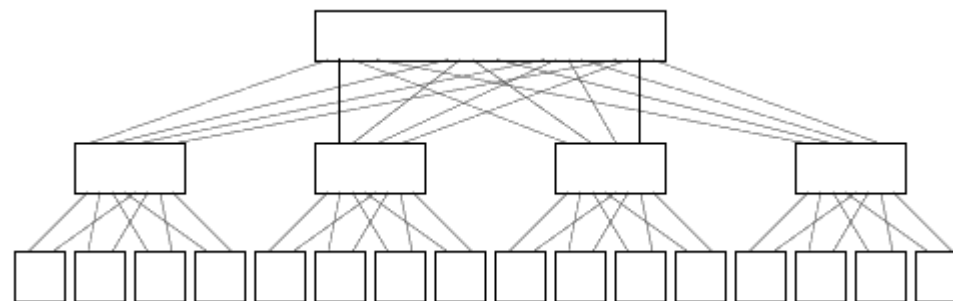
**d:** 1 (koncové uzly), 2 (kořenový), 3 (ostatní)  
**D:**  $2\log((N+1)/2)$   
**B:** 1  
**C:** N-1

**Binary tree**

# Statické sítě



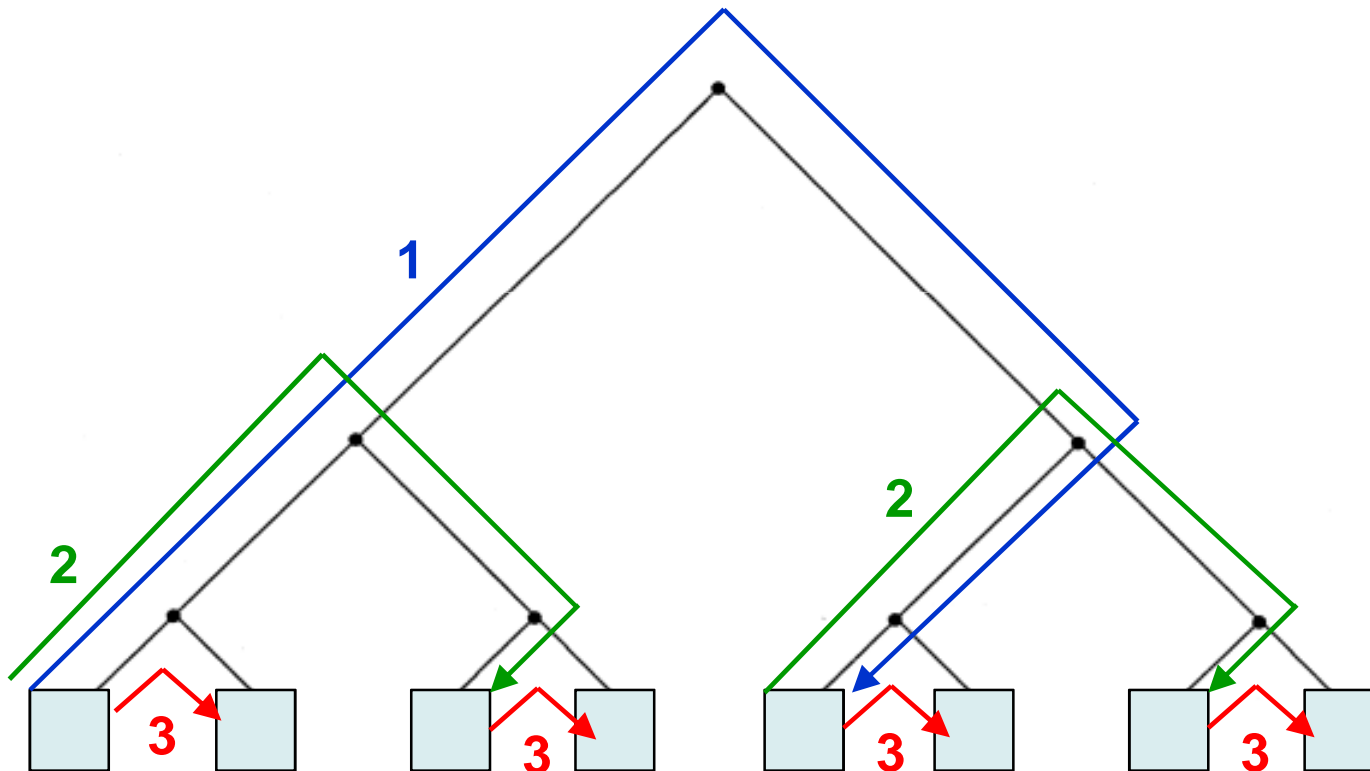
**Binary fat tree**



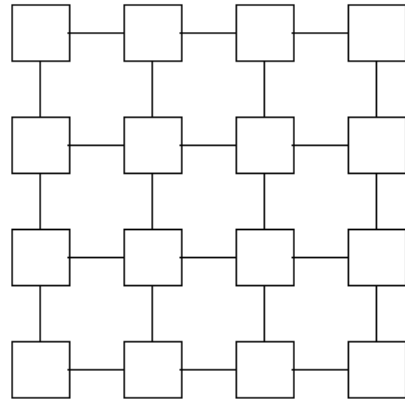
**4-ary fat tree**

## Nepřímý Binární strom – příklad komunikace

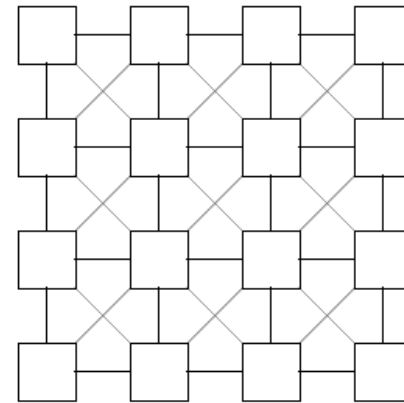
- V nepřímém stromu jsou výpočetní uzly v až listech...
- One-to-all broadcast a All-to-one redukce



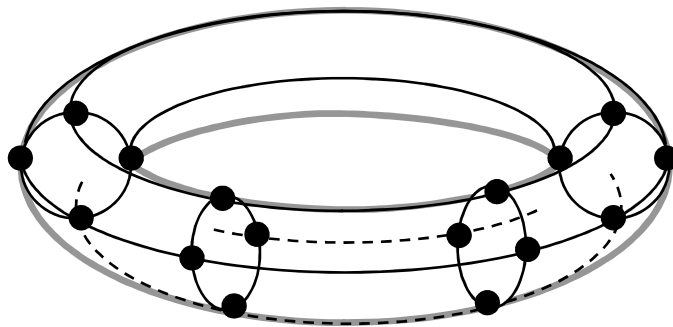
# Statické sítě



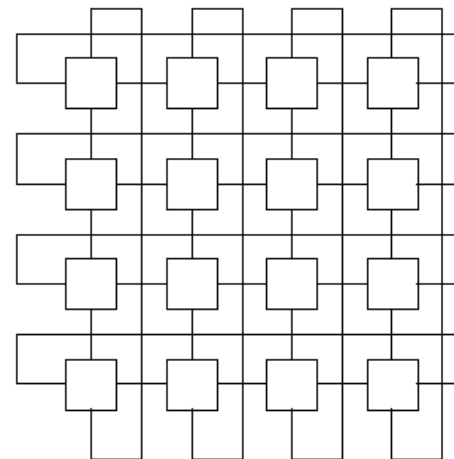
**4-connected 2D mesh**



**8-connected 2D mesh**



**2D Torus (anuloid)**

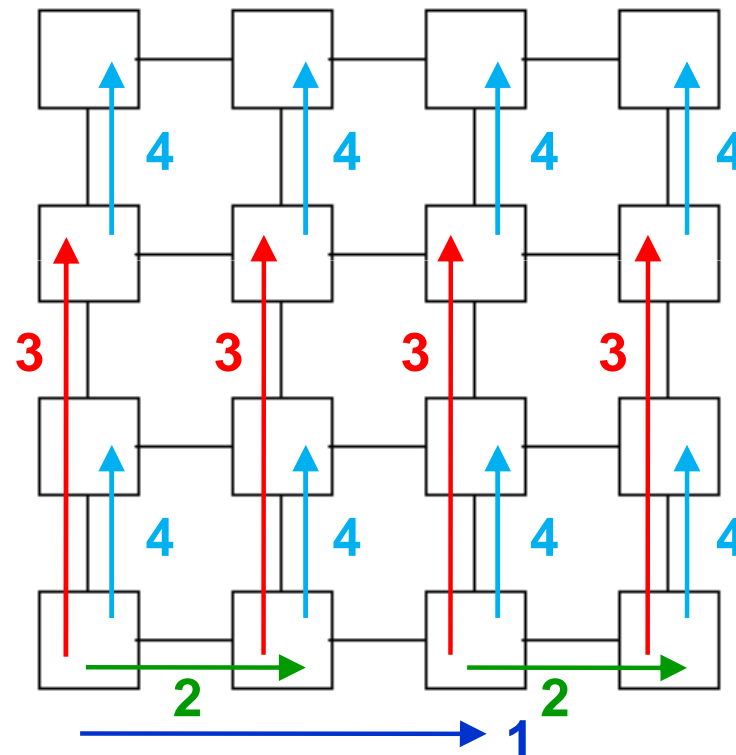


d: 4  
D:  $N^{1/2}$   
B:  $2N^{1/2}$   
C:  $2N$

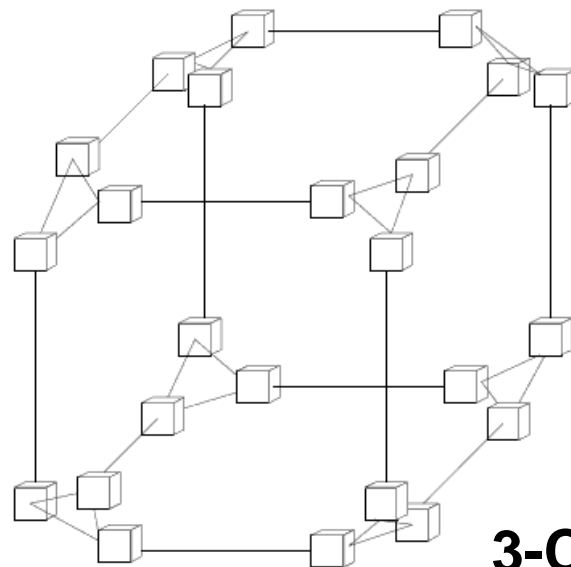
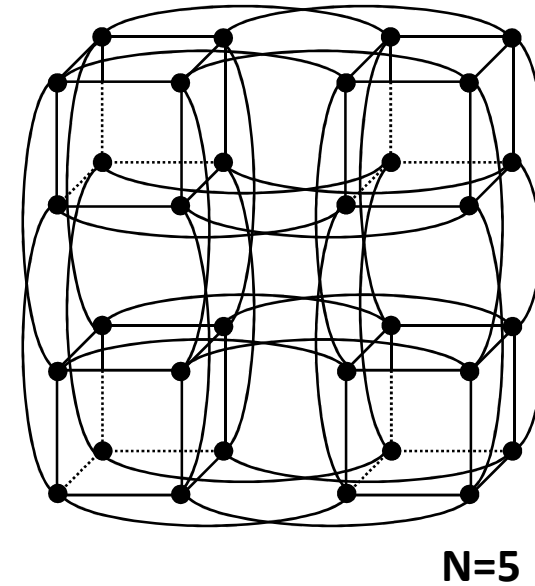
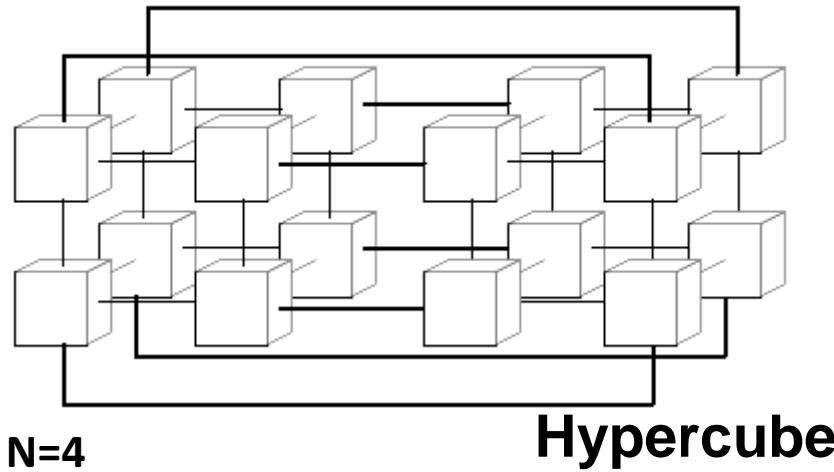


## Mesh – příklad komunikace

- One-to-all broadcast a All-to-one redukce

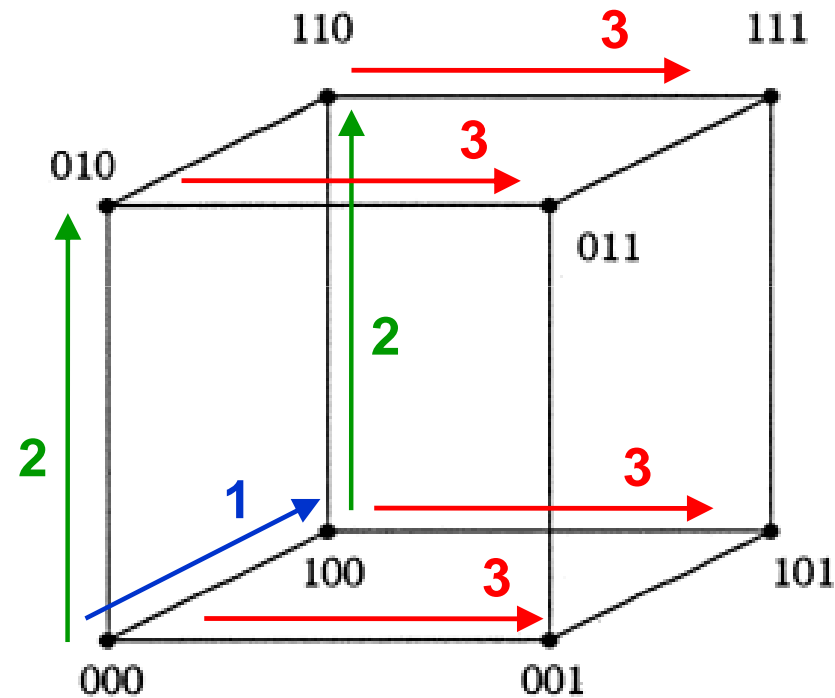


# Statické sítě

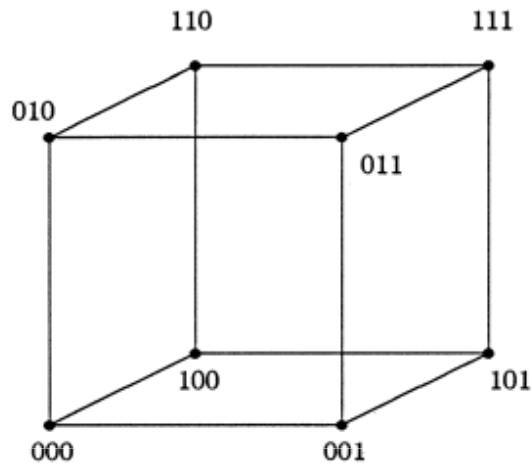


## Hypercube – příklad komunikace

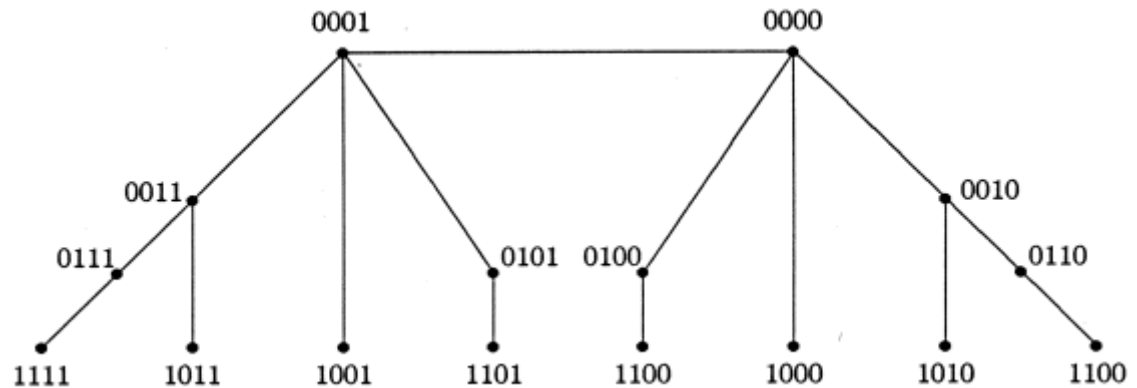
- One-to-all broadcast a All-to-one redukce



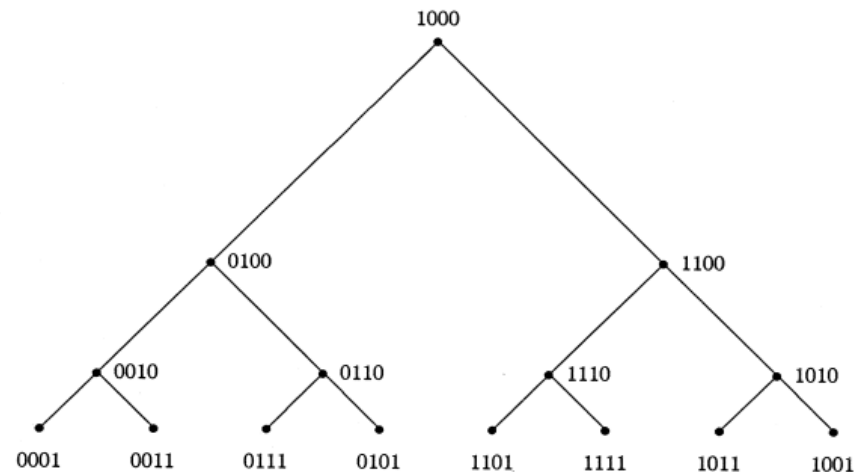
# Hypercube – vztah k jiným sítím



Stromová struktura 4-dimenzionální krychle:

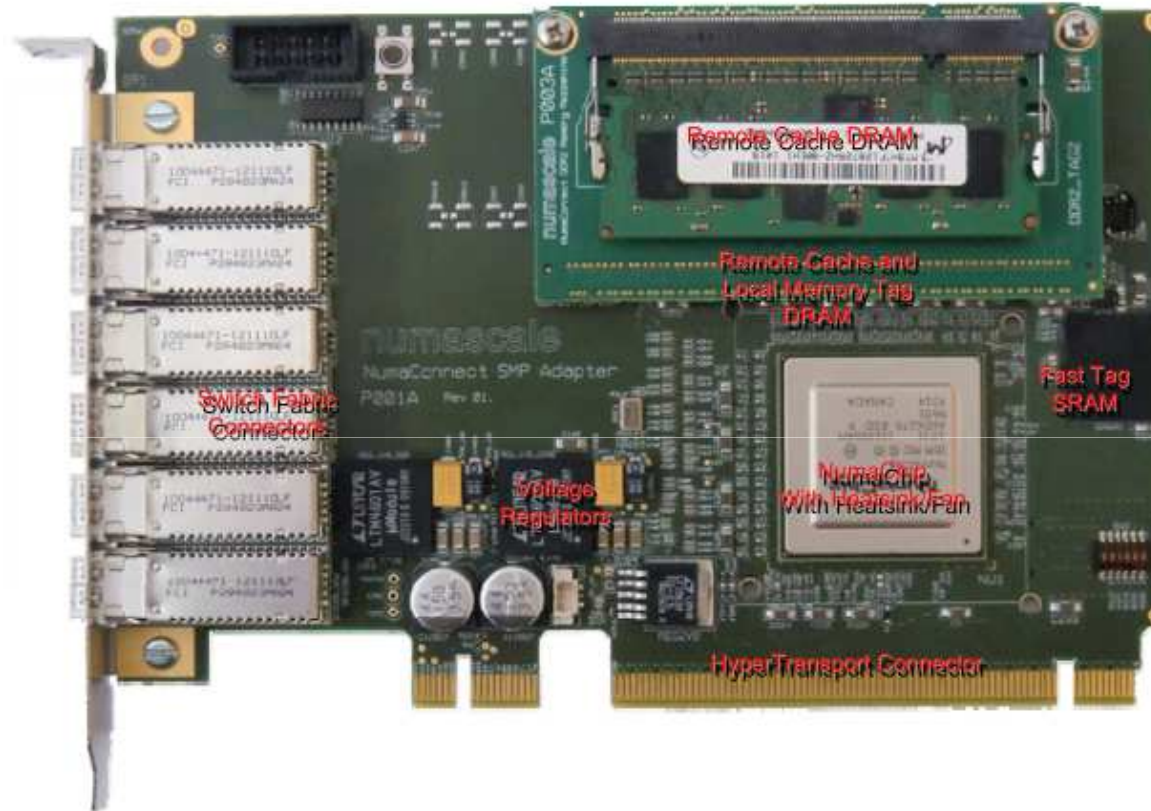


Pro  $n \geq 3$  je možné mapovat  $n$ -úrovňový binární strom do podgrafu získaného z hyperkrychle při odstranění některého z uzlů. Hrana pak reprezentuje buď hranu hyperkrychle nebo cestu délky 2.



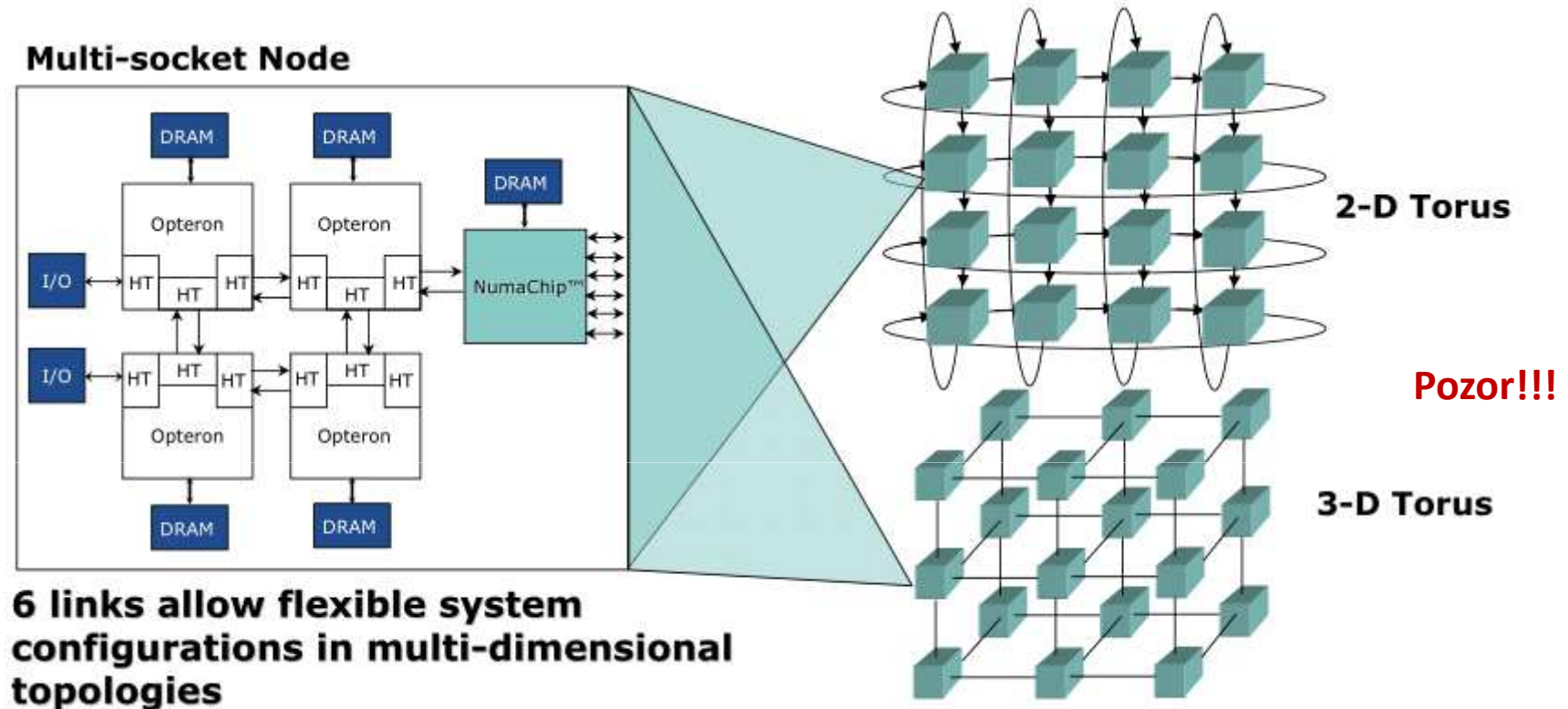
## Diskuze – Tím to však nekončí...

numascale



- Scalable directory based cache coherence in hardware
- Support for 4096 Nodes
- 4GB Cache 8 GB Tag (supports 240 GB Local Node RAM)

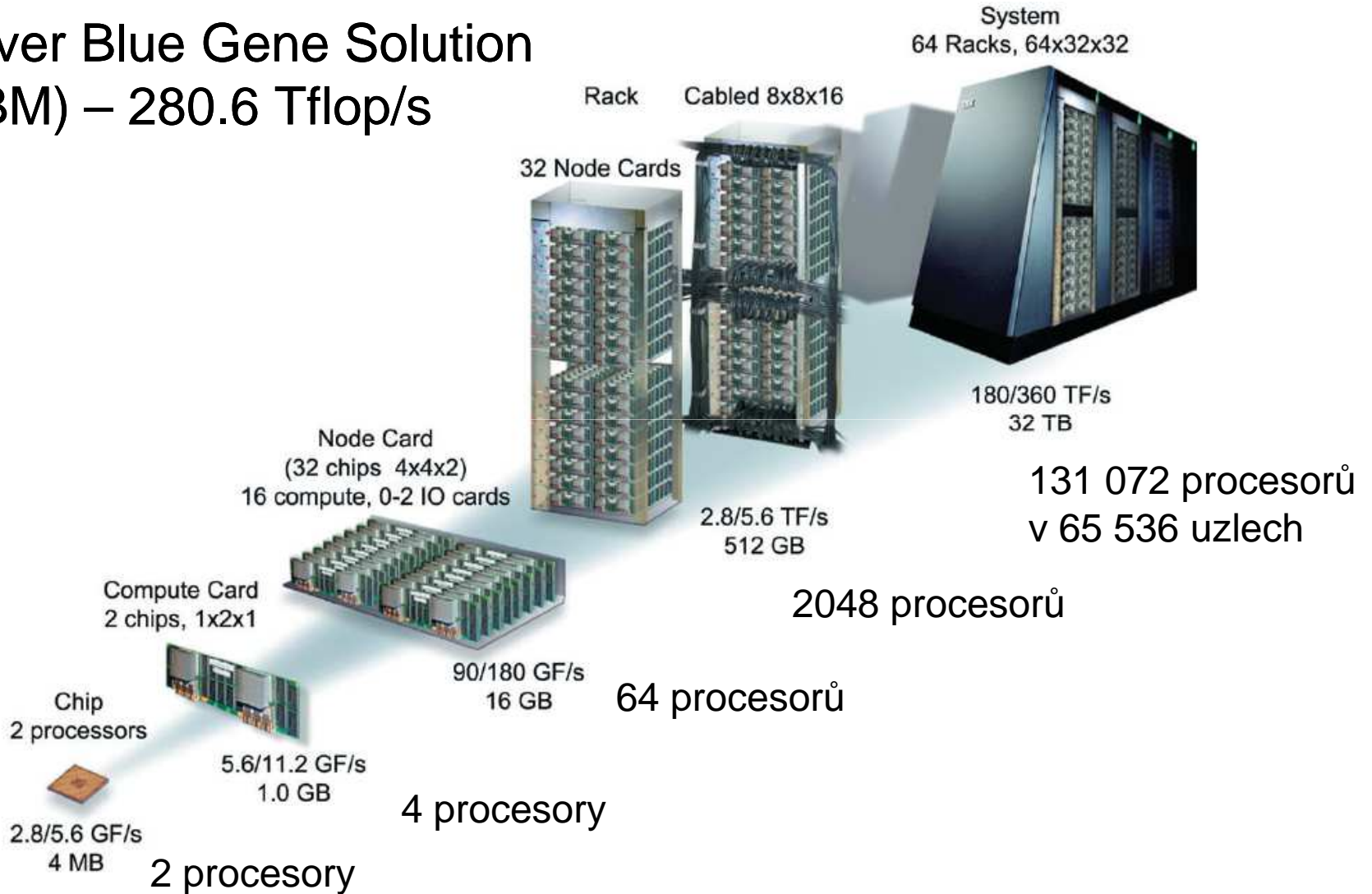
## Diskuze – Tím to však nekončí...



- Use any Programming Model Available for the Node on the whole System (OpenMP, MPI, Threads, ...)
- NO Application Changes Required!

# Motivace - Blue Gene Solution

## eServer Blue Gene Solution (IBM) – 280.6 Tflop/s



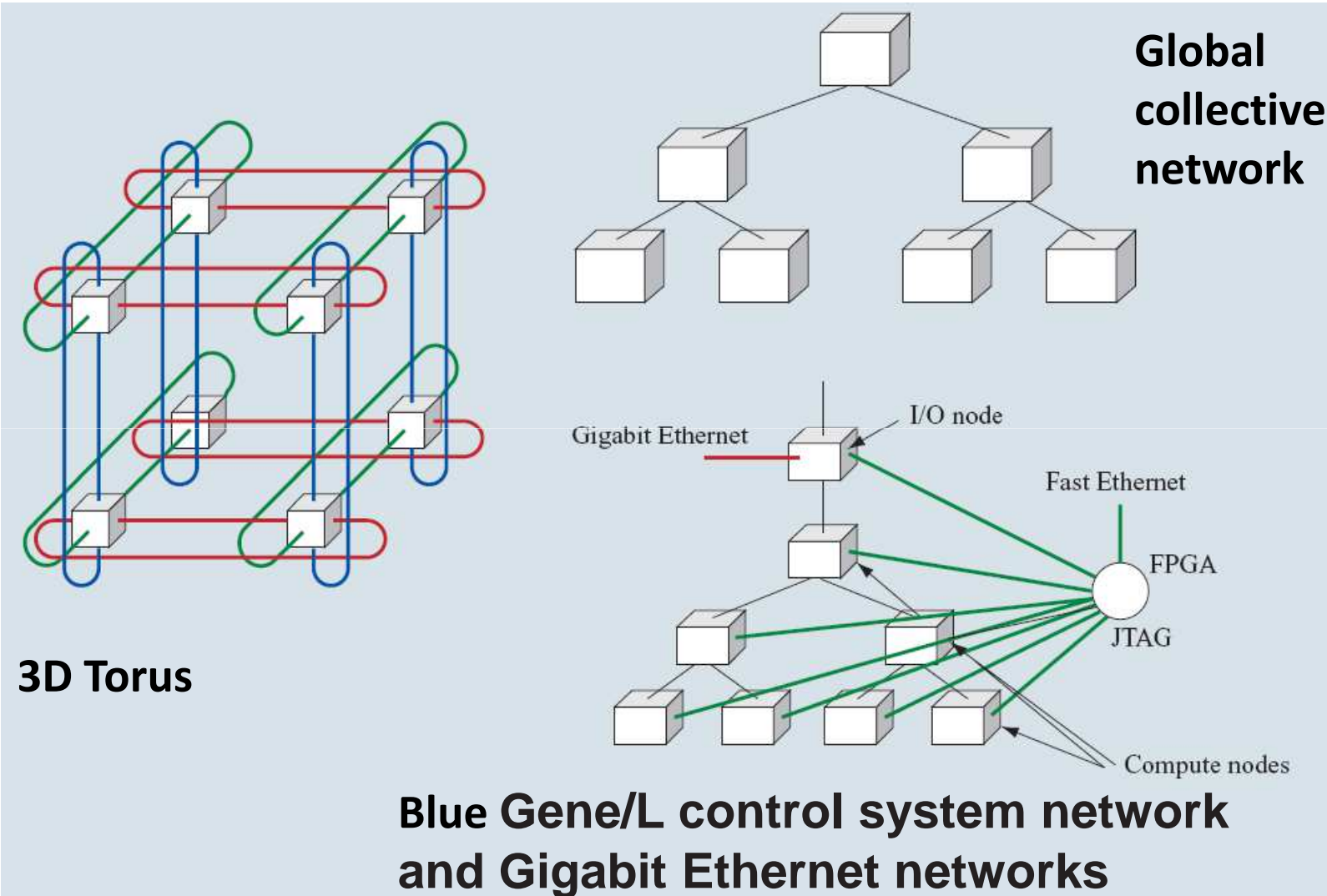
## Motivace - Blue Gene Solution

Využívá 5 sítí pro propojení uzlů

- **3D torus pro point-to-point komunikaci mezi uzly (175 MBps v každém směru),**
- **globální hromadní síť pro hromadné operace (Global collective network) – 350 MBps, 1.5 $\mu$ s latence (data z jednoho uzlů mohou být rozeslána všem ostatním uzlům – broadcast, nebo jenom některým),**
- **global barrier and interrupt network,**
- **control network (system boot, debug, monitoring stavu teploty, ventilátorů,...)**
- **gigabit Ethernet network pro řízení a I/O operace**



# Motivace - Blue Gene Solution



# Topologie

## Topologie propojovacích sítí

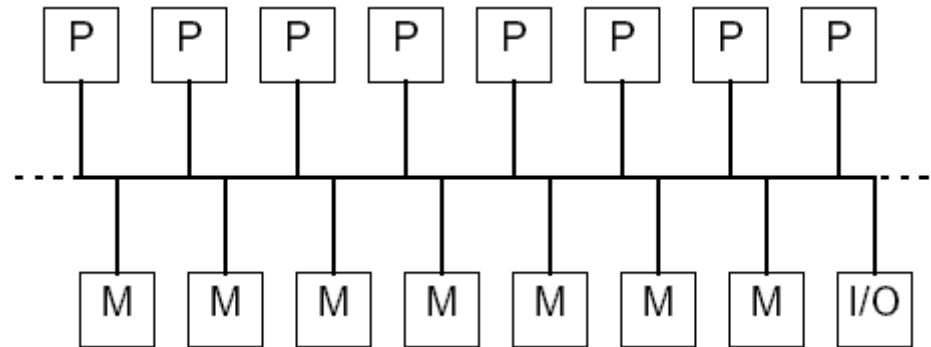
### Statické sítě

- lineární pole (linear array)
- kruh (ring)
- Chordalový kruh (Chordal ring)
- binary tree
- fat tree
- 2D, 3D mesh
- 2D, 3D torus
- hypercube
- Cube Connected Cycles (CCC)

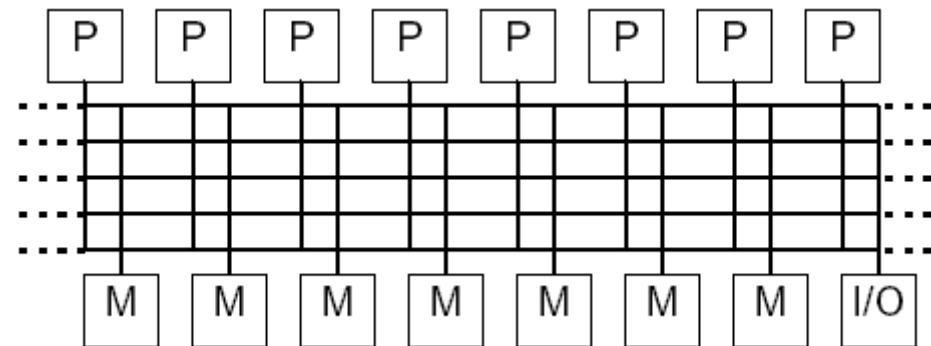
### Dynamické sítě

- Bus network
- jednostupňové sítě (křížové prepínače - crossbars)
- vícestupňové sítě (omega, Banyan, Cantor, Clos,...)

## Dynamické sítě

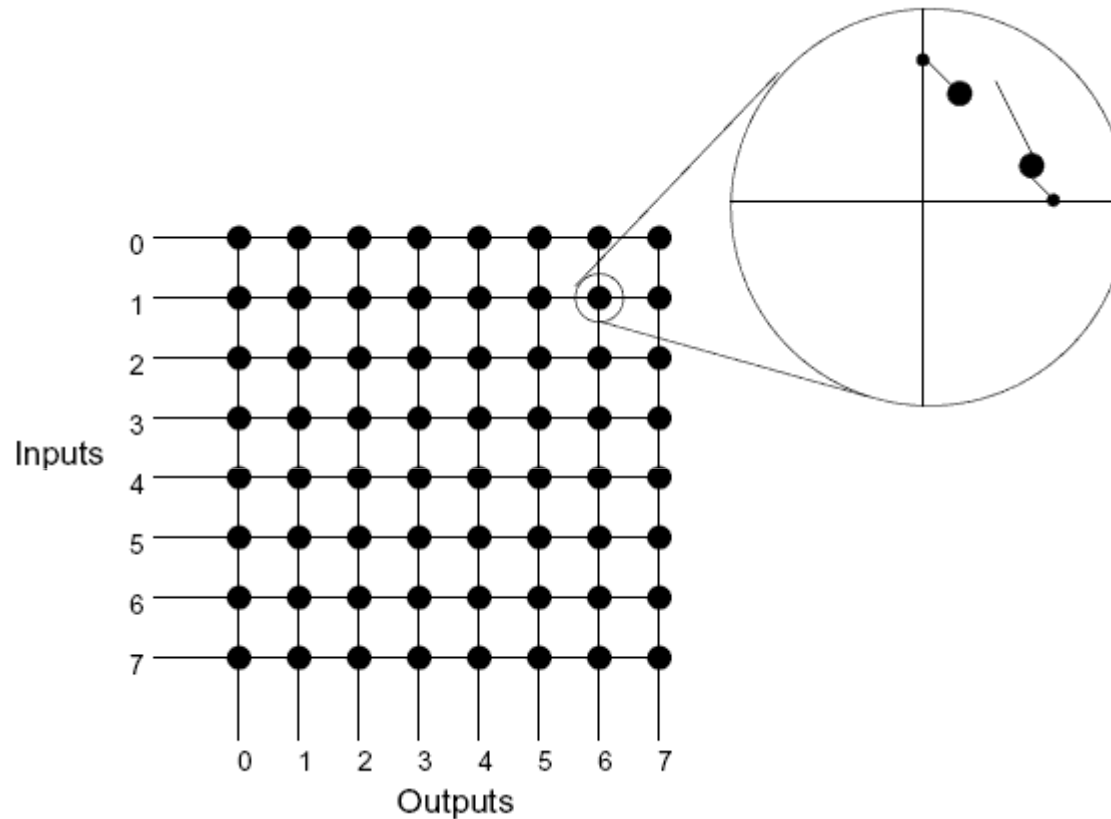


**Bus-based multiprocessor system**



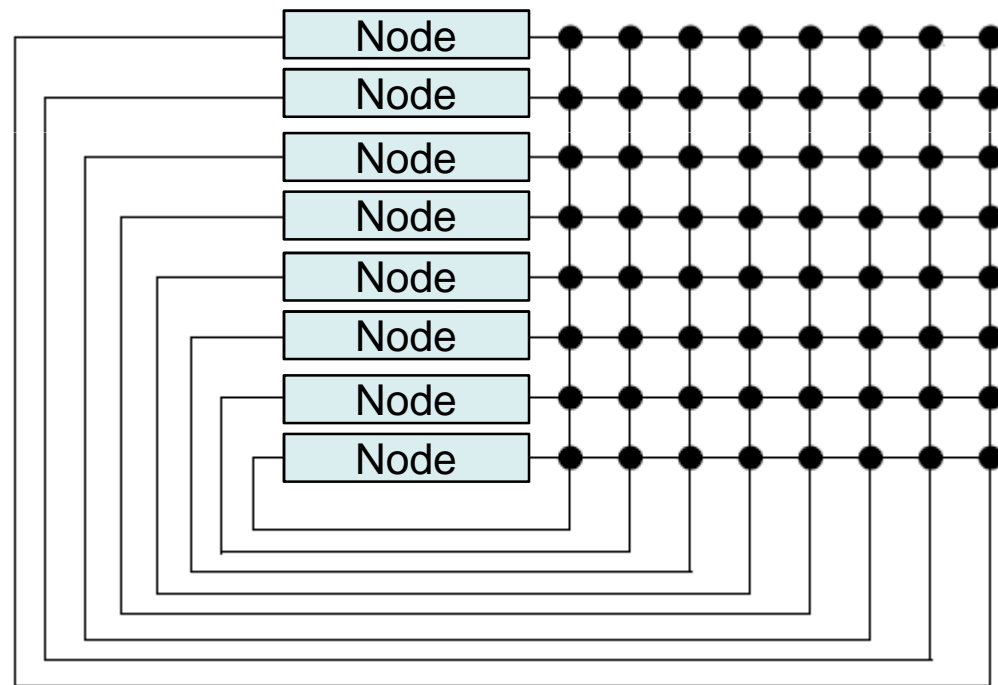
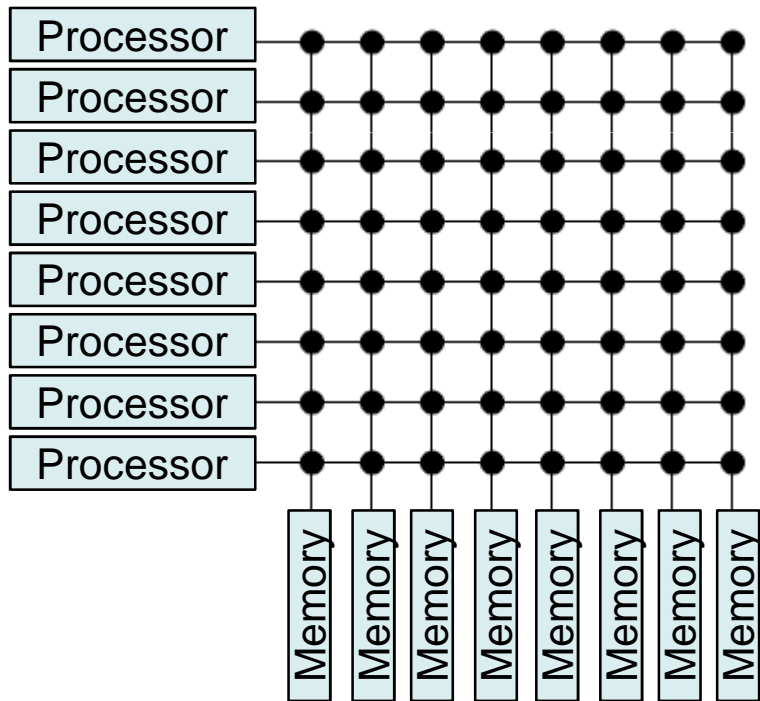
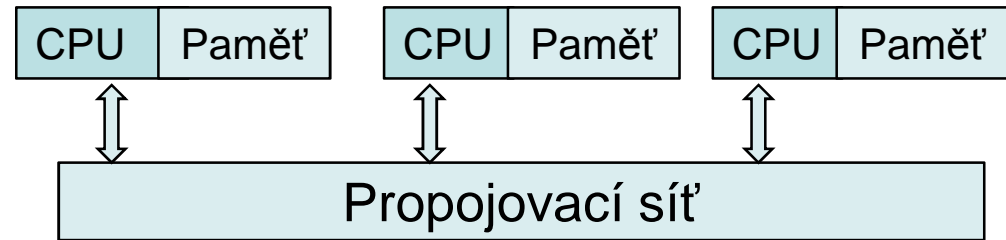
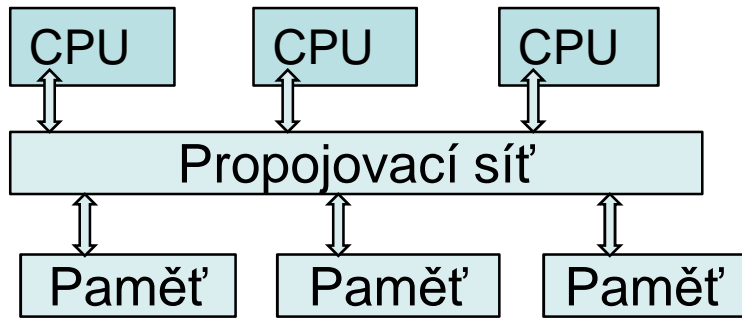
**Multiple bus architecture**

# Dynamické sítě – Single stage network

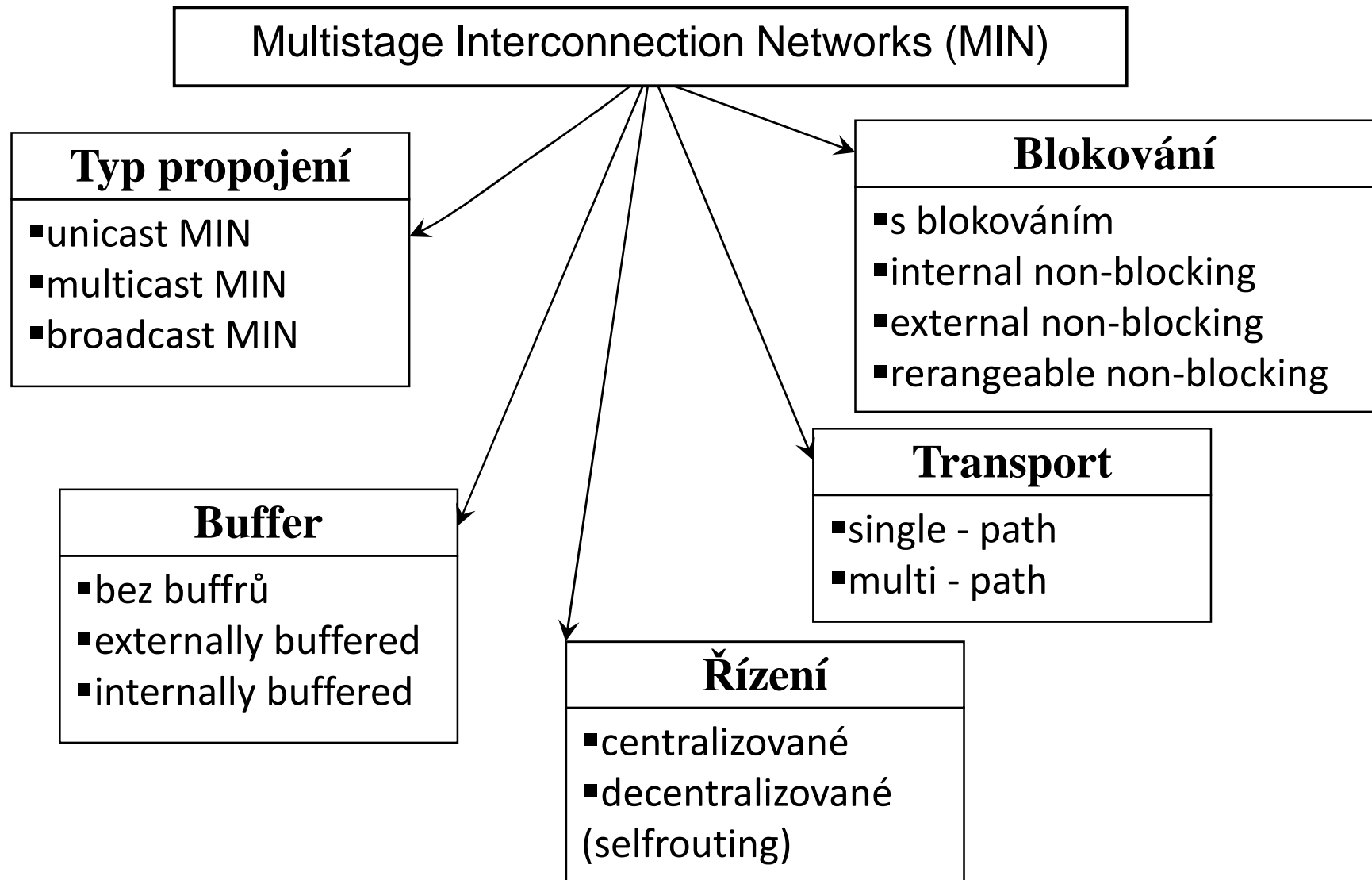


**Crossbar – Křížový přepínač – Cena  $N^2$**

# Dynamické sítě – Single stage network



# Dynamické sítě – Multistage interconnection networks



## Dynamické sítě – Multistagestage interconnection networks

- **Unicast MIN** (one-to-one, point-to-point) - je síť, která propojuje v téže době právě jeden vstup s právě jedním výstupem.
- **Multicast MIN** (multipoint) - propojí v téže době jeden vstup sítě s více výstupy sítě
- **Broadcast MIN** - propojuje jeden vstup sítě se všemi výstupy sítě v téže době. Broadcast MIN je speciálním případem Multicast MIN.
  
- **Jednocestní MIN** (Single-Path) -všechny pakety v rámci jednoho virtuálního spojení používají stejnou cestu.
- **Vícecestní MIN** (Multi-Path) -pakety v rámci jednoho virtuálního spojení používají několik cest. Pakety jsou náhodně distribuovány. Vnitřní provoz ve spojovací síti se stane nezávislým na vnějším provozu (tj na poměru služeb, které jsou propojovány v síti). Taková síť vyžaduje opatření pro řazení buněk u východu ze sítě.

## Dynamické sítě – Multistage interconnection networks

- **MIN s centralizovaným řízením** -řízení je prováděno centrálním procesorem. Procesor určuje výběr cesty v síti, pokud jsou dány požadavky na spojení.
- **MIN s decentralizovaným řízením** -řízení je distribuováno na spojovací elementy (samosměrování - self routing). V tomto případě potřebujeme po vstupu do sítě pro každý paket přídatnou informaci, tag, na základě které se v každém elementu určuje kombinace vstupu a výstupu.



## Dynamické sítě – Multistage interconnection networks

- **Sítě s blokováním** - jsou sítě kdy se konkrétní vstup nemůže propojit s konkrétním výstupem, přesto že výstup není obsazen jiným spojením. Blokování je spojeno se ztrátou informace, nebo její zpožděním. Pro zajištění kvality spojování je žádoucí blokování vyloučit, nebo jej minimalizovat na specifikovanou úroveň.
- **Sítě bez vnitřního blokování** (Internal non-blocking) - jsou sítě, které zajišťují propojení libovolného vstupu s libovolným výstupem bez zrušení, nebo rekonfigurace jiné vnitřní cesty v síti. Může však nastat případ, že i když síť je bez vnitřního blokování, dva, nebo více vstupů může chtít použít stejný výstup v téže době - vzniká konflikt.
- **Sítě bez vnějšího blokování** (External non-blocking) - jsou sítě schopné propojení libovolného vstupu s libovolným výstupem v každém případě. V síti je nutná v případě konfliktu funkce uložení informace, tj vyrovnávací paměť.
- **Rekonfigurovatelné** (preuspořadatelné) **sítě bez blokování** (Rearrangeable non-blocking) - jsou zvláštním případem blokovacích sítí - ovšem dokáží vždy realizovat spojení z libovolného vstupu na libovolný výstup. V případě konfliktu dochází k restrukturalizaci (přeuspořádání) stávajících cest a k vytvoření nové konfigurace spojení.

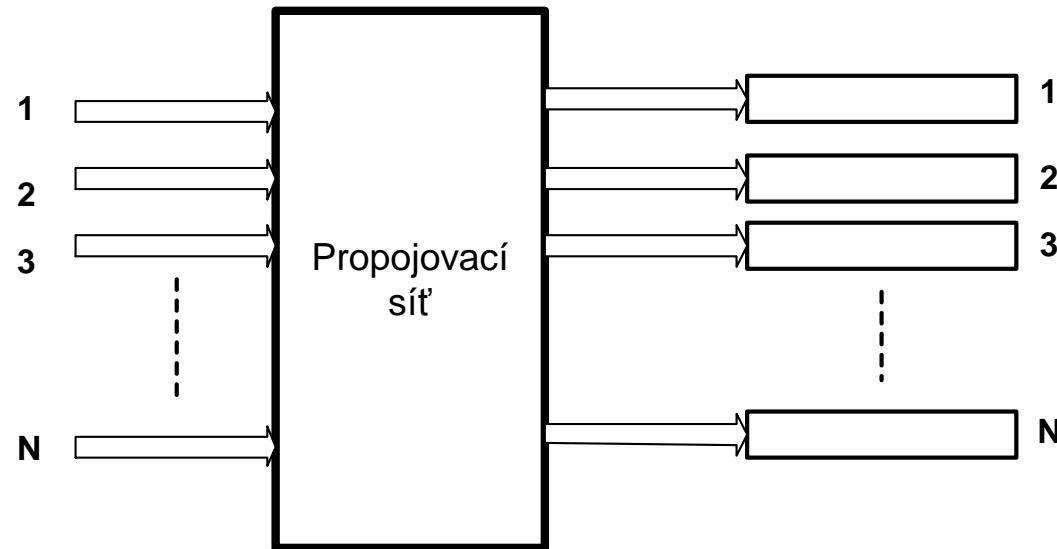
## Dynamické sítě – Multistage interconnection networks

- **MIN s vnější vyrovnávací pamětí** (externally Buffered) - vyrovnávací paměti jsou umístěny na vstupech sítě, na výstupech sítě, nebo kombinovaně.
- **MIN s vnitřní vyrovnávací pamětí** (Internally Buffered) - vyrovnávací paměti jsou v jednotlivých spojovacích elementech, čili uvnitř sítě. I v spojovacím elementu může být vyrovnávací paměť na vstupu (Input Buffering), na výstupu (Output Buffering), nebo uprostřed elementu (Central Buffering)

### **Hlavní způsoby řazení (ukládání) paketů:**

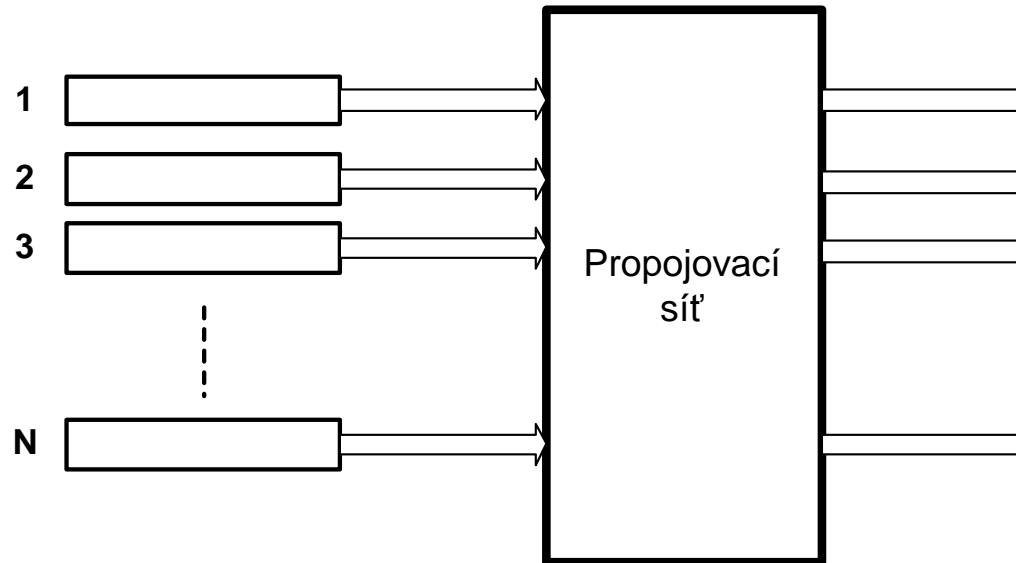
- výstupní řazení (OQ)
- vstupní řazení (IQ)
- kombinované vstupně-výstupní řazení (CIOQ)
- centralizované sdílené řazení (CSQ)
- virtuální výstupní řazení (VOQ)

## Výstupní řazení paketů: (IQ - output queuing)



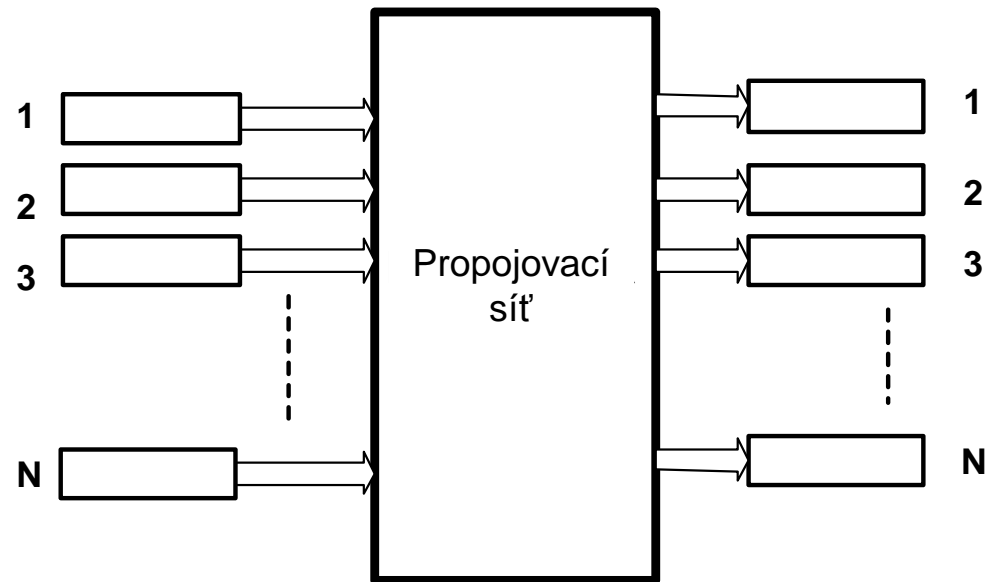
Když paket dorazí na vstupní port, je okamžitě (po průchodu sítí) uložen do bufferu, který je na příslušném výstupním portu. Protože pakety určené pro stejný výstupní port mohou přijít současně z mnoha vstupních portů, výstupní buffer potřebuje radit pakety mnohem větší rychlostí než vstupní port. V nejhorším možném případě to může být až  $N$  krát rychleji (kde  $N$  je počet vstupních portů), a to tehdy pokud pakety ze všech vstupních portů jsou určeny pro jeden určitý výstupní port. Rychlost, kterou můžeme přistupovat do výstupního bufferu je však limitována.

## Vstupní řazení paketů: (IQ - output queuing)



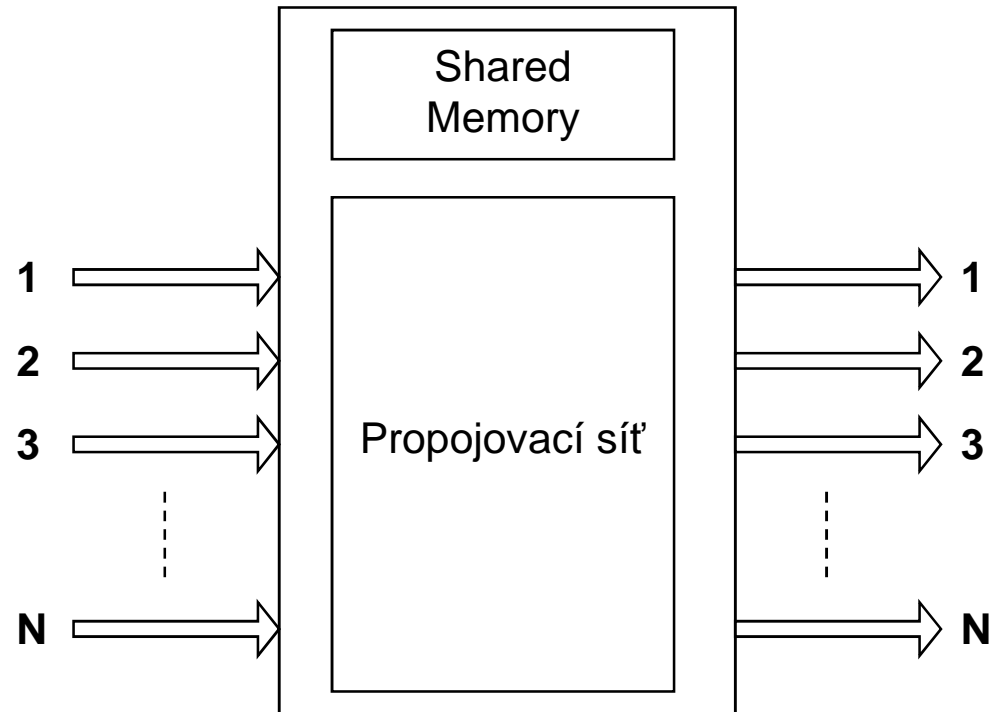
Vstupní řazení nemá limity jaké má například výstupní řazení nebo centralizované sdílené řazení. V této architektuře má každý vstupní port FIFO zásobník buněk, a pouze první buňka ve frontě je způsobilá pro přenos v průběhu daného časového úseku. Nevýhodou FIFO řazení je, že když buňka na čele fronty se zablokuje, zablokují se všechny buňky za ní a zabrání se tak jejich přenosu, dokonce i když je výstupní port volný. Toto se nazývá Head-of-line blokování. Matematickou analýzou a počítačovou simulací bylo ukázáno, že HOL blokování omezuje propustnost každého vstupního portu na maximálně 58.6 procent a to při náhodné hustotě provozu, a tato hodnota je ještě mnohem nižší při velmi hustém provozu.

## Kombinované vstupně-výstupní řazení paketů: (COQ - combined input-output queuing)



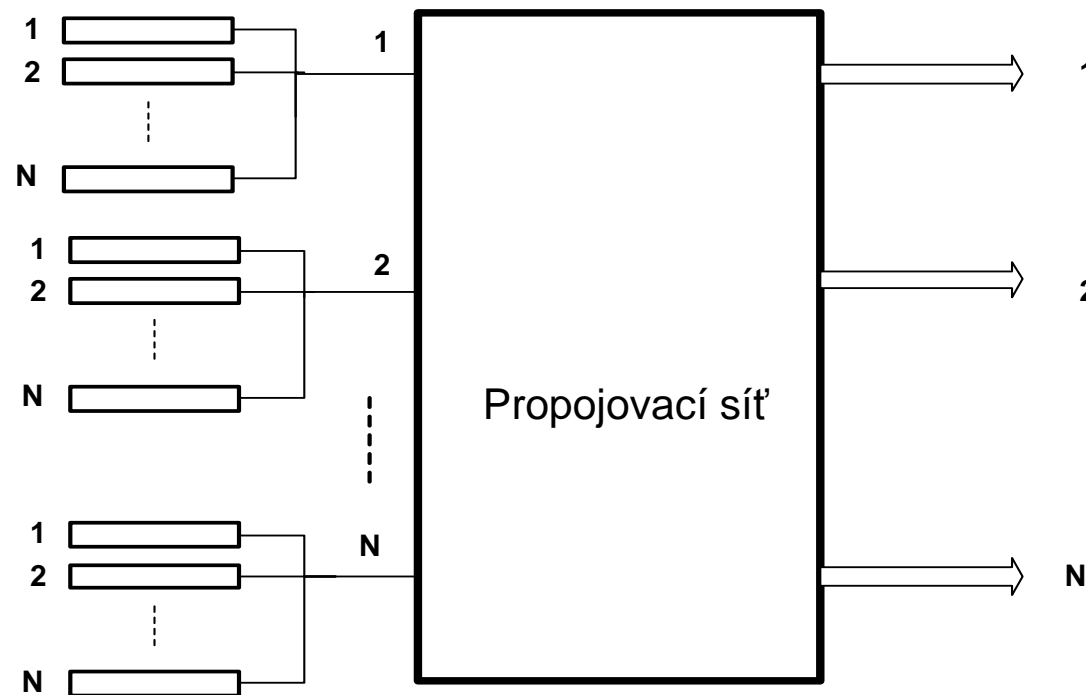
Toto schéma řazení je kombinací vstupního a výstupního řazení. Je to dobrý kompromis mezi výkonem a možností rozšiřování OQ a IQ přepínačů. Pro vstupně řazené přepínače, nanejvýš jeden paket může být doručen na výstupní port v jednom časovém úseku. Pro výstup frontované přepínače, až N paketů může být doručených na výstupní port za jeden časový úsek. Použitím CIOQ, namísto těchto dvou extrémních možností, můžeme si zvolit kompromis mezi nimi.

## Centralizované sdílené řazení paketů: (Centralized shared queuing)



Máme buffer sdílený všemi výstupními porty přepínače, na který se můžeme dívat jako na sdílenou paměťovou jednotku, a který má N souběžných přístupů k zápisu pro N vstupních portů, a N souběžných přístupů na čtení pro N výstupních portů. Protože pakety určené pro stejný výstupní port mohou přicházet současně z více vstupních portů, výstupní port musí být schopen načítat mnohem rychleji než vstupní port dokáže zapisovat údaje.

## Virtuální výstupní řazení paketů: (VO - virtual output queuing)



Toto schéma řazení zvládá Head-of-line blokování a zároveň si udržuje svou výhodu rozšiřitelnosti. Při této metodě každý vstupní port udržuje izolovanou frontu pro každý výstupní port. Klíčovým faktorem pro dosažení vysokého výkonu použitím VOQ přepínačů je plánovací algoritmus, který je zodpovědný za výběr paketů, které by měly být přeneseny v každé časové jednotce ze vstupních portů na výstupní. Několik takových algoritmů již bylo navrženo, jako např. PIM (parallel iterativní matching - paralelní iterační párování), Islip a RPA. Bylo ukázáno, že s méně než čtyřmi opakováními výše uvedeného řídicího algoritmu PIM, propustnost přepínače přesáhne 99 procent.

# Základní typy vícecestupňových sítí

## Unicast síť s blokováním jednocestné (Single-Path)

- Unicast spojovací síť provádějí **one-to-one** spojení. Jednocestné spojovací síť (Single-Path) mají právě jen jednu cestu a tím jen jednu možnost spojení mezi libovolným vstupem a libovolným výstupem.
  - Baseline síť
  - Banyan síť
  - Delta síť
  - Omega síť

## Unicast síť s blokováním vícecestné (Multi-Path)

- Vícecestné síť zajišťují alternativní cesty mezi vstupy a výstupy, přičemž zachovávají samosměřovací vlastnosti sítě a jen minimálně komplikují její časovou složitost. Zlepší se spolehlivost i propustnost sítě.
  - Banyan síť s dělenou zátěží
  - Data Manipulator - DM
  - Augmented Data Manipulator - ADM
  - Reverse Augmented Data Manipulator - IADM
  - gama síť



# Základní typy víceúrovňových sítí

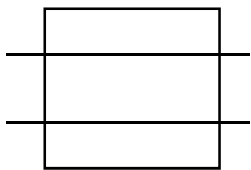
## Unicast síť bez blokování

- Všechny sítě s blokováním vyžadují opatření na potlačení blokování. Nejčastěji je to umístění vyrovnávacích pamětí do sítě.
- Problém blokování je možné řešit již při návrhu sítě, vytvořením sítě bez blokování. Rozlišujeme dvě možnosti. Buď jsou sítě topologicky bez blokování, tj jejich architektura minimalizuje pravděpodobnost blokování, nebo jsou s řízením bez blokování, tj mechanismus řízení sítě odstraňuje blokování. Druhý případ nazýváme i rekonfigurovatelné sítě.
  - Benešova síť
  - Paralelní baseline síť
  - Closova síť
  - Batcherova síť
  - Batcher-banyan síť

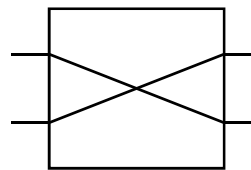
# Základní typy vícestupňových sítí

## Multicast síť

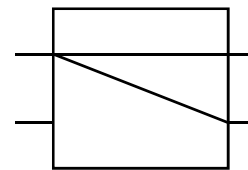
- Multicast vícestupňové spojovací sítě jsou obvykle sítě s  $N$  vstupy a  $N$  výstupy, přičemž **libovolná skupina vstupů se může propojit s libovolnou skupinou výstupů**. Každý vstupní port může být spojen s více než jedním výstupem, ale každý výstupní port je obvykle spojen s nejvíce jedním vstupním portem.
- Multicast sítě mohou uskutečnit  $N^N$  různých spojení, mají tedy větší výkon než unicast sítě, které uskutečňují one-to-one permutaci vstupů na výstupy a realizují maximálně  $N!$  různých spojení.
- V podstatě jako multicast síť může fungovat každá unicast síť pokud její spojovací elementy dokáží propojit své vstupy na několik výstupů.



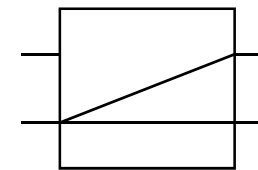
Straight-trough



Exchange

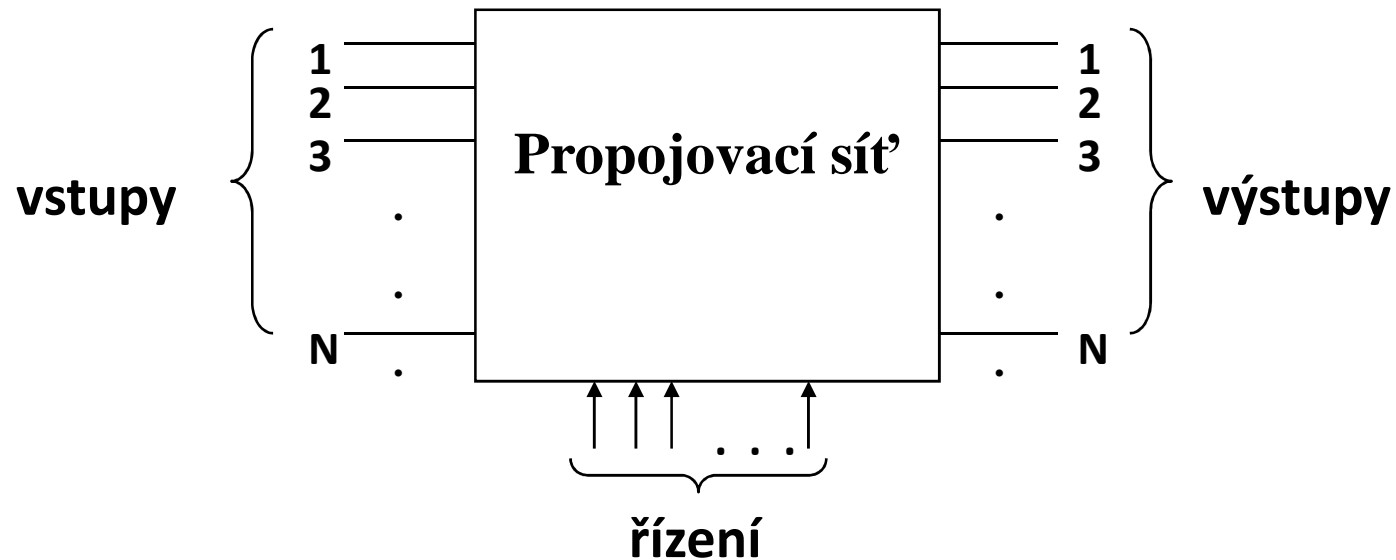


Upper-broadcast



Lower-broadcast

## Blokové schéma propojovací sítě



- Propojovací síť - PS  $[N \times N]$  je prostředek, který umožňuje propojit libovolný ze svých  $N$  vstupů ( $X_0, X_1, \dots, X_{N-1}$ ) s libovolným ze svých  $N$  výstupů ( $Y_0, Y_1, \dots, Y_{N-1}$ ).
- Základní význam mají PS, které realizují pouze navzájem jednoznačné přiřazení, tzn. s jedním vstupem může být spojen pouze jediný výstup -> **permutační síť**

# Základní permutace

Základní permutace:

- **permutace dokonalého promíchání** ( $\sigma$ ) - perfect shuffle
- **motýlková permutace** ( $\beta$ ) - Butterfly
- **permutace reverzace** ( $\rho$ )
- **permutace výměny** ( $E$ ) - exchange

- Permutace dokonalého promíchání:

$$\sigma(j) = \left( 2j + \left\lfloor \frac{2j}{N} \right\rfloor \right) \bmod N$$

- Obecně je operace promíchání definovaná ve tvaru:

$$\sigma(j, K) = \left( Kj + \left\lfloor \frac{Kj}{N} \right\rfloor \right) \bmod N$$

- kde  $K$  je "počet rukou" potřebných na promíchání, pro které platí  $K = 2^k$ ,  $k = 0, 1, \dots, n-1$ .

## Permutace dokonalého promíchání

- Pokud máme číslo  $x$  v binární reprezentaci, pak permutace dokonalého promíchání odpovídá cyklickému posuvu binární reprezentace o jeden bit vlevo:

$$\sigma(\mathbf{x}) = \sigma([x_n \ x_{n-1} \ x_{n-2} \ \dots \ x_2 \ x_1]) = [x_{n-1} \ x_{n-2} \ \dots \ x_1 \ x_n]$$

- Permutace dokonalého  $k$ -tého podpromíchání  $\sigma_{(k)}(\mathbf{x})$  - je definována pro  $1 \leq k \leq n$  následovně:

$$\sigma_{(k)}(\mathbf{x}) = \sigma_{(k)}([x_n \ x_{n-1} \ x_{n-2} \ \dots \ x_2 \ x_1]) = [x_n \ \dots \ x_{k+1} \ x_{k-1} \ \dots \ x_1 \ x_k]$$

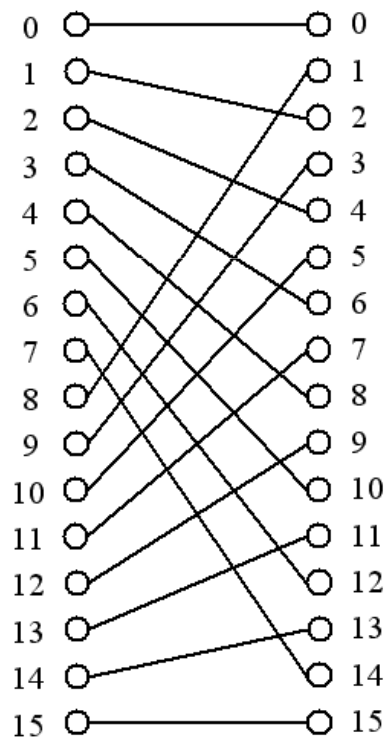
- Takže platí:

$$\sigma_{(1)}(\mathbf{x}) \equiv \mathbf{x}, \quad \sigma_{(n)}(\mathbf{x}) \equiv \sigma(\mathbf{x})$$

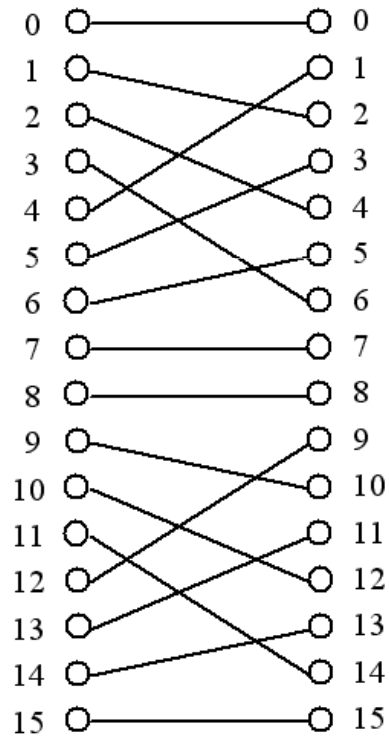
# Permutace dokonalého promíchání

Například pro  $N=16$ ,  $n=\log_2 N=4$ :

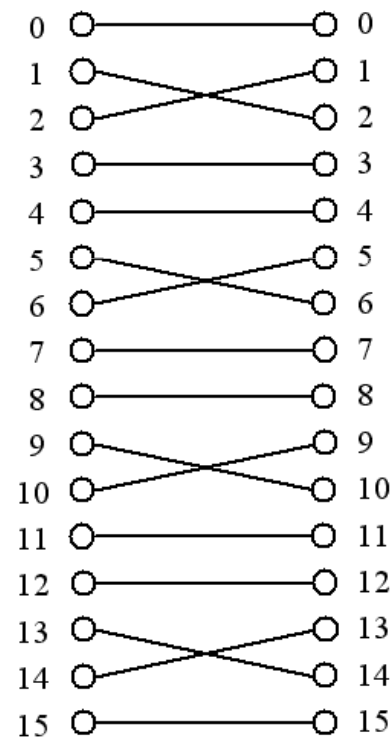
- $\sigma_{(4)}(0) = \sigma_{(4)}([0000]) = [0000] = 0$
- $\sigma_{(4)}(1) = \sigma_{(4)}([0001]) = [0010] = 2$
- $\sigma_{(4)}(2) = \sigma_{(4)}([0010]) = [0100] = 4$  atd.



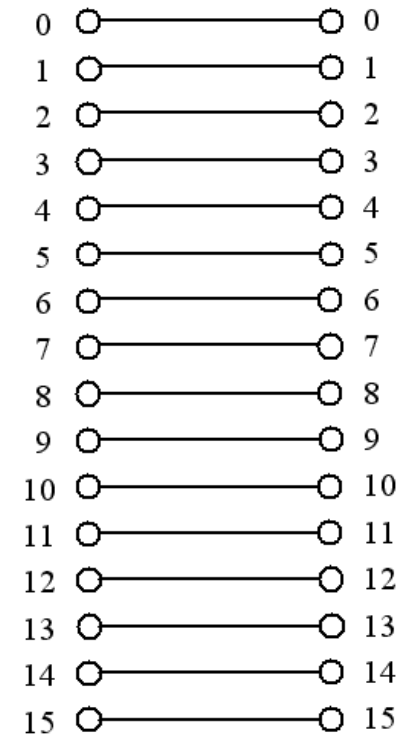
$\sigma_{(4)}(x) \equiv \sigma(x)$



$\sigma_{(3)}(x)$



$\sigma_{(2)}(x)$

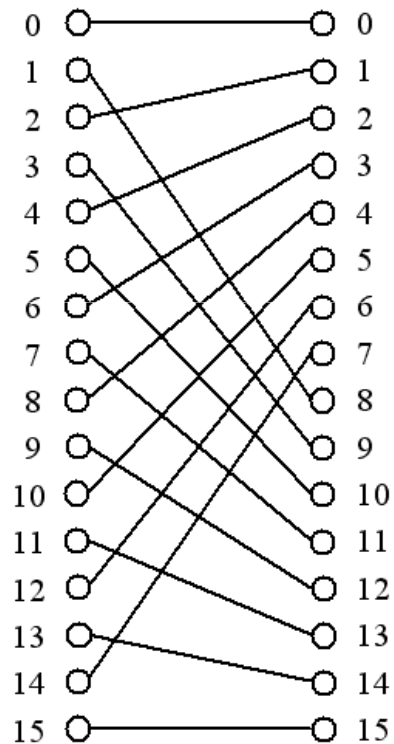


$\sigma_{(1)}(x) \equiv x$

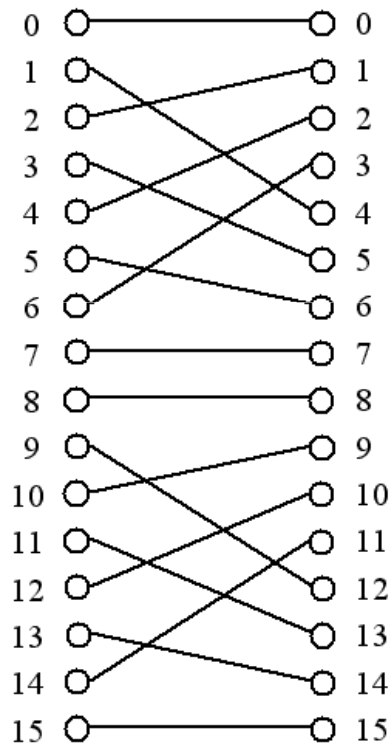
# Permutace dokonalého promíchání

**Inverzní dokonalé promíchání** odpovídá cyklickému posuvu binární reprezentace o jeden bit vpravo:

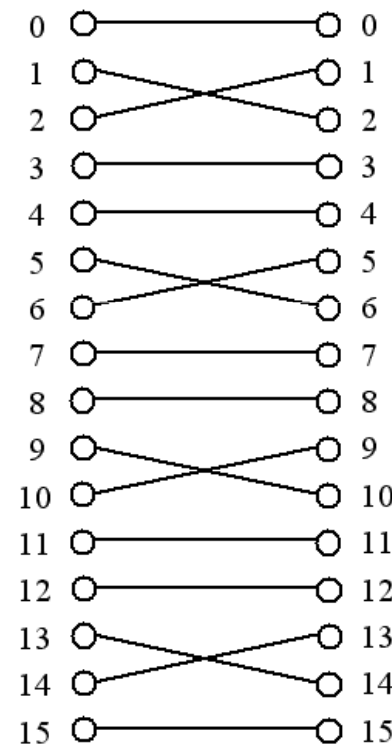
$$\sigma^{-1}(\mathbf{x}) = \sigma^{-1}([x_n \ x_{n-1} \ x_{n-2} \ \dots \ x_2 \ x_1]) = [x_1 \ x_n \ x_{n-1} \ x_{n-2} \ \dots \ x_2]$$



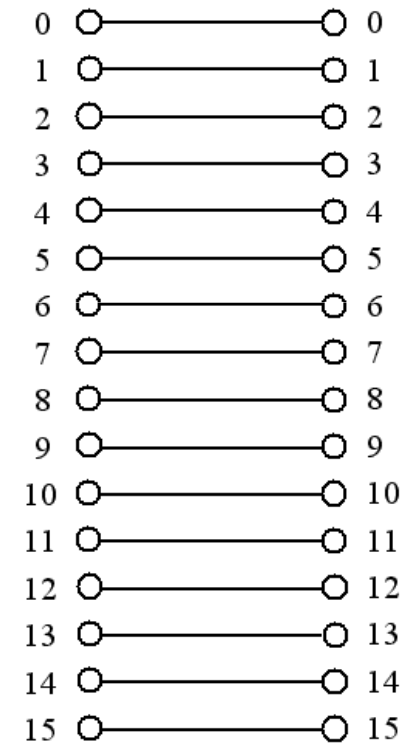
$$\sigma_{(4)}^{-1}(\mathbf{x}) \equiv \sigma^{-1}(\mathbf{x})$$



$$\sigma_{(3)}^{-1}(\mathbf{x})$$



$$\sigma_{(2)}^{-1}(\mathbf{x})$$



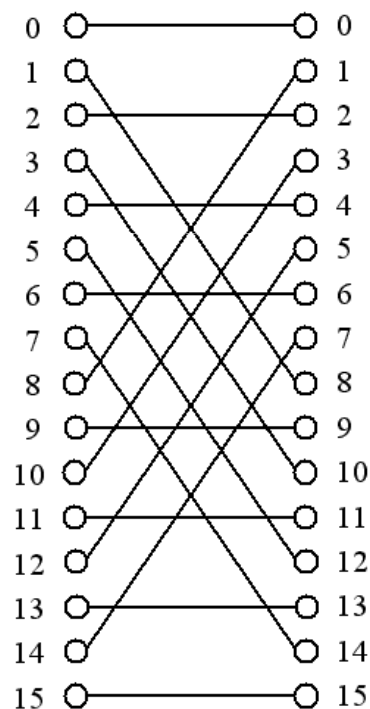
$$\sigma_{(1)}^{-1}(\mathbf{x}) \equiv \mathbf{x}$$

# Motýlková permutace

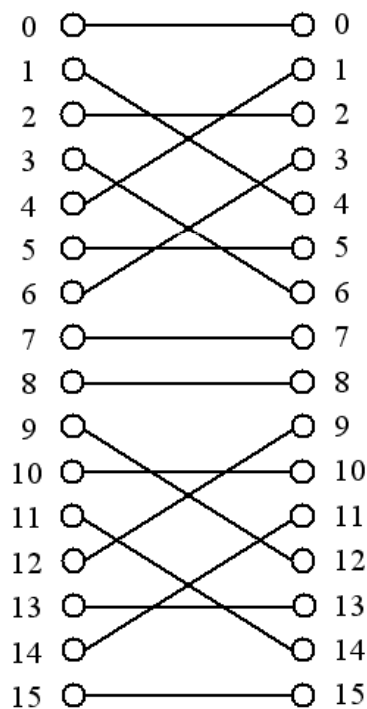
## Motýlková permutace (Butterfly):

k-ta motýlková permutace  $\beta_{(k)}(x)$  pro  $1 \leq k \leq n$ ,  $n = \log_2 N$  se dostane tak, že se vzájemně vymění první a k-tý bit:

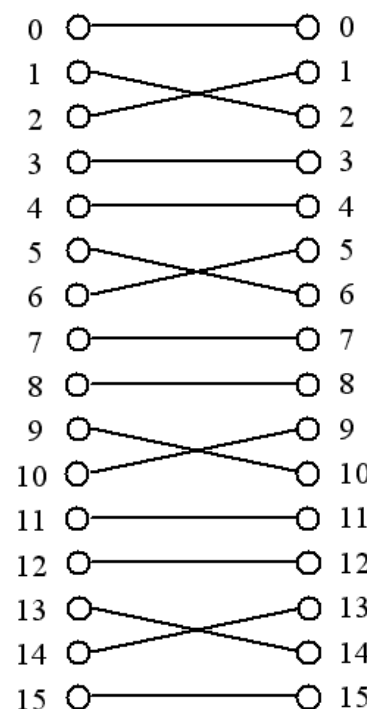
$$\beta_{(k)}(\mathbf{x}) = \beta_{(k)}([x_n \ x_{n-1} \ x_{n-2} \ \dots \ x_2 \ x_1]) = [x_n \ x_{n-1} \ \dots \ x_{k+1} \ x_1 \ x_{k-1} \ \dots \ x_2 \ x_k]$$



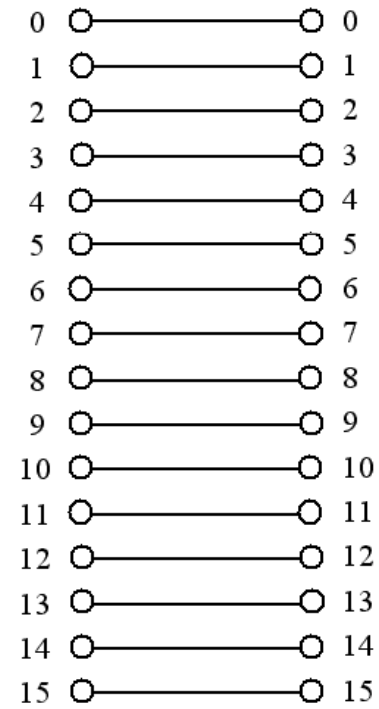
$$\beta_{(4)}(\mathbf{x}) \equiv \beta_{(4)}^{-1}(\mathbf{x})$$



$$\beta_{(3)}(\mathbf{x}) \equiv \beta_{(3)}^{-1}(\mathbf{x})$$



$$\beta_{(2)}(\mathbf{x}) \equiv \beta_{(2)}^{-1}(\mathbf{x})$$



$$\beta_{(1)}(\mathbf{x}) \equiv \beta_{(1)}^{-1}(\mathbf{x})$$

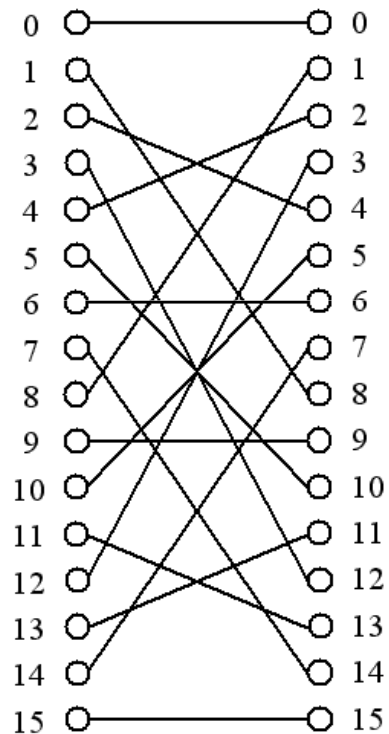


# Permutace reverzace

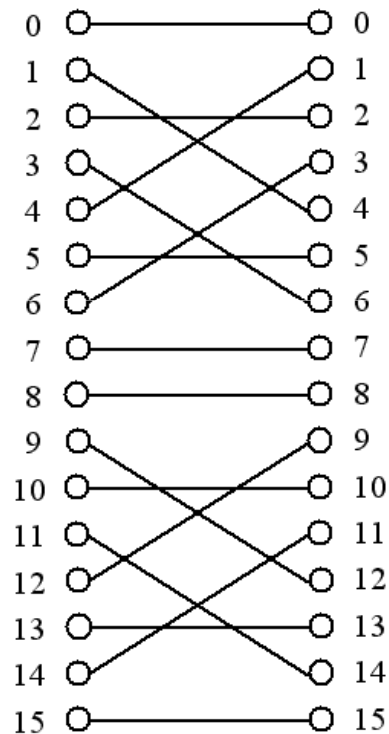
## Permutace reverzace:

je dána vzájemnou výměnou bitů vyšších a nižších (zrcadlení):

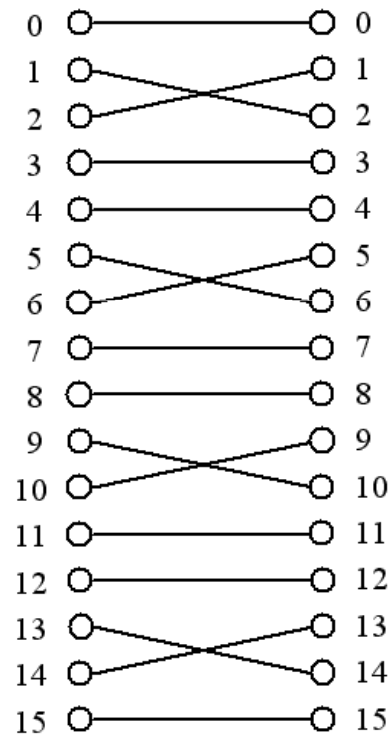
$$\rho(\mathbf{x}) = \rho([x_n \ x_{n-1} \ \dots \ x_2 \ x_1]) = [x_1 \ x_2 \ \dots \ x_{n-1} \ x_n]$$



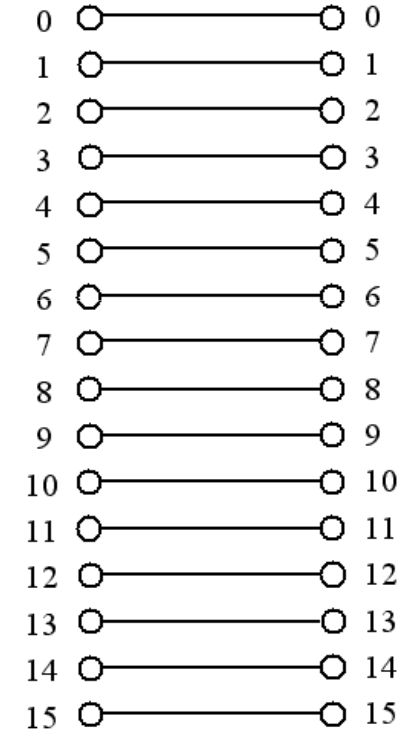
$\rho = \rho_4$



$\rho_3$



$\rho_2$



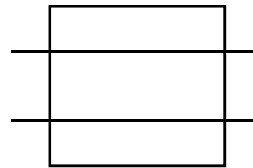
$\rho_1 = X$

## Permutace výměny

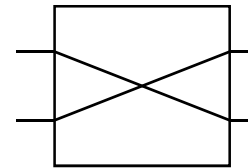
### Permutace výměny:

Definujme  $\hat{x}$  vztahem  $\hat{x} = [x_n \ x_{n-1} \ \dots \ x_2 \ \bar{x}_1]$ , kde  $\bar{x}_1$  znamená inverzi prvního (nejnižšího) bitu. Pak množina výměnných permutace je:

$$\begin{array}{l} \text{bud' } e(x) = x \quad \text{a} \quad e(\hat{x}) = \hat{x}, \\ \text{nebo } e(x) = \hat{x} \quad \text{a} \quad e(\hat{x}) = x. \end{array}$$



**Straight-trough**  
(na přímo)

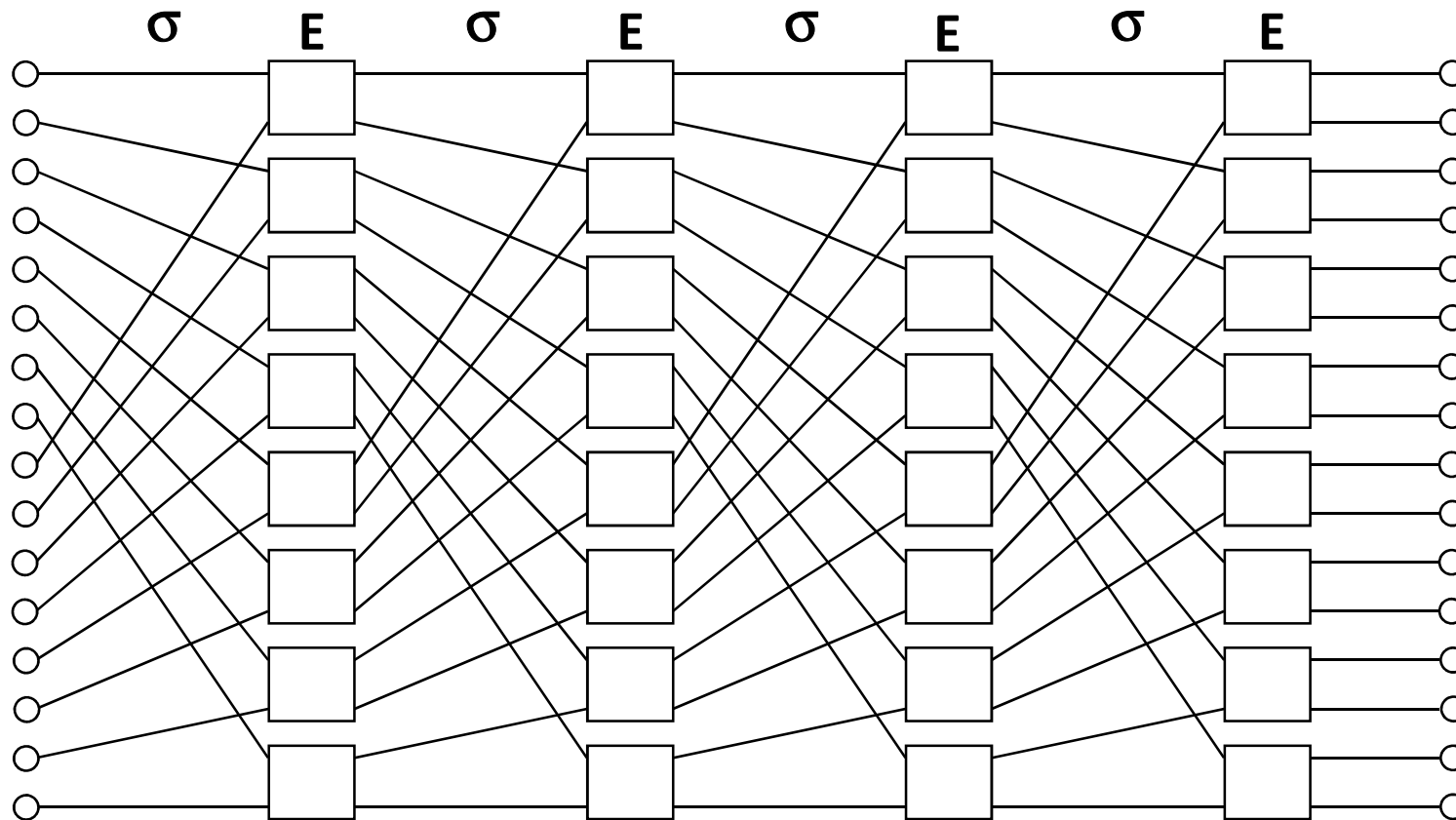


**Exchange**  
(výměna)

# Omega síť

$$\Omega_N = \sigma E \sigma E \dots \sigma E = (\sigma E)^n \quad n = \log_2 N$$

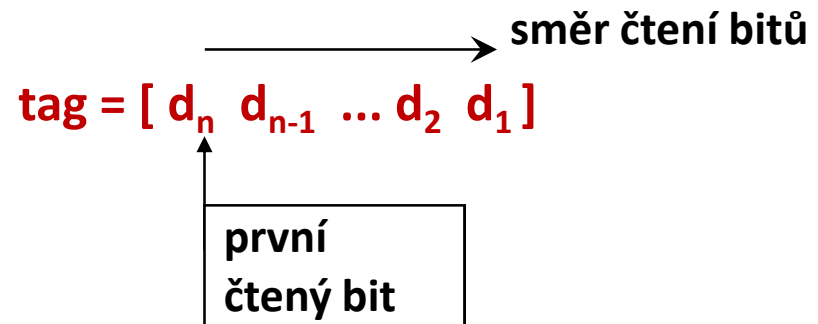
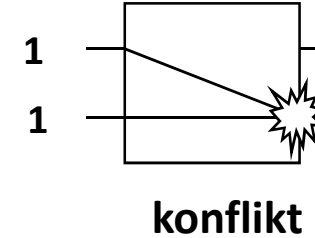
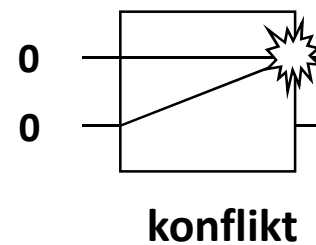
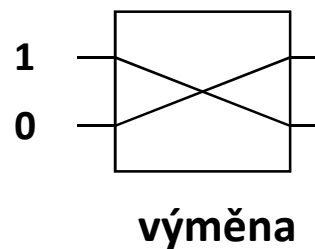
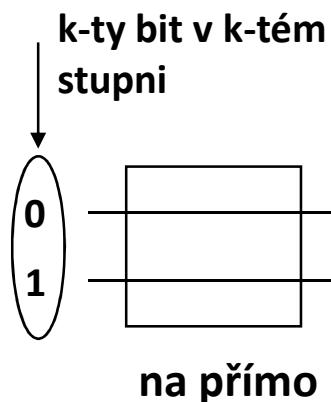
- Napr. pre  $N=16$ ,  $n=4$ :  $\Omega_{16} = (\sigma E)^4$



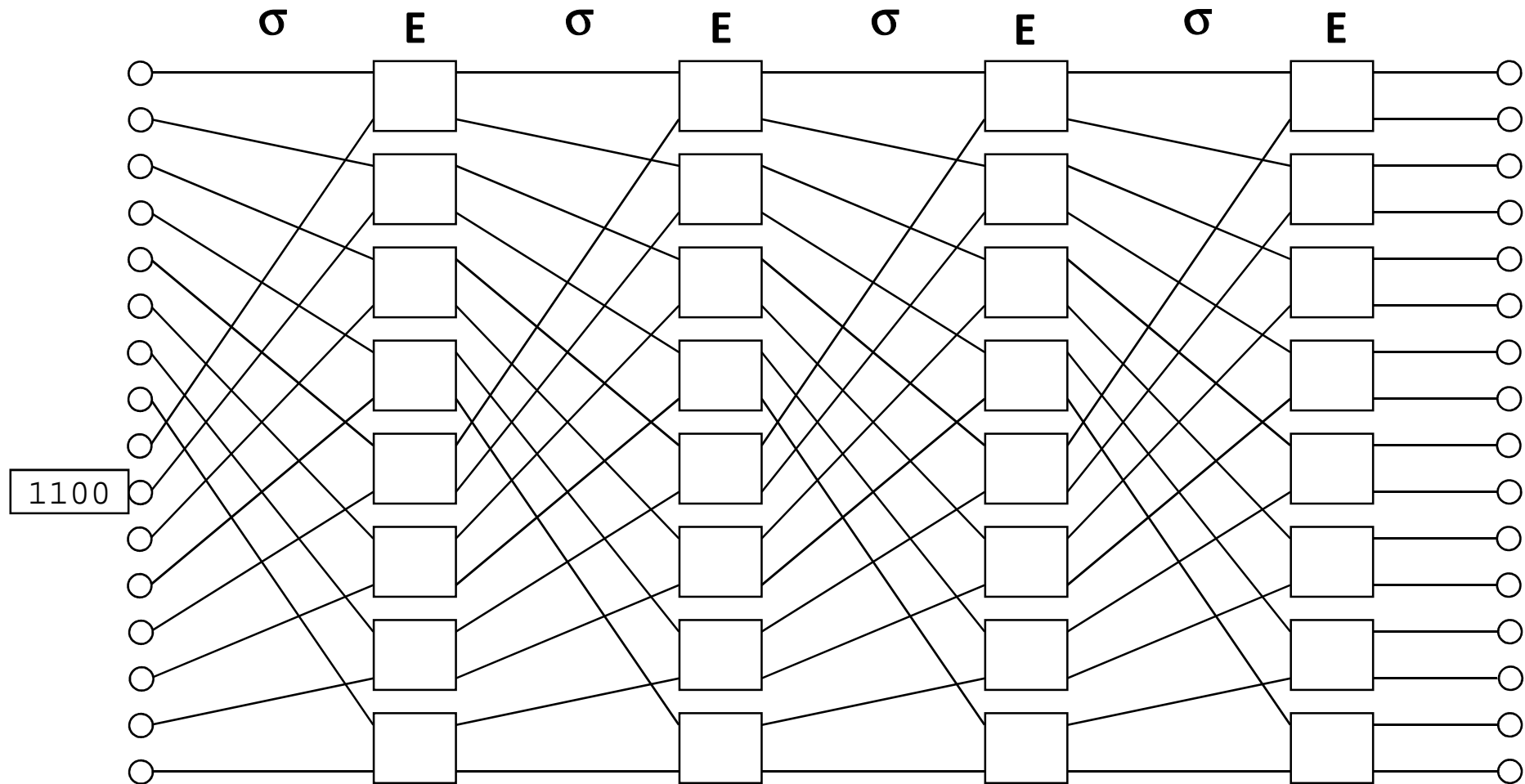
# Omega síť

## Samosměrovací algoritmus:

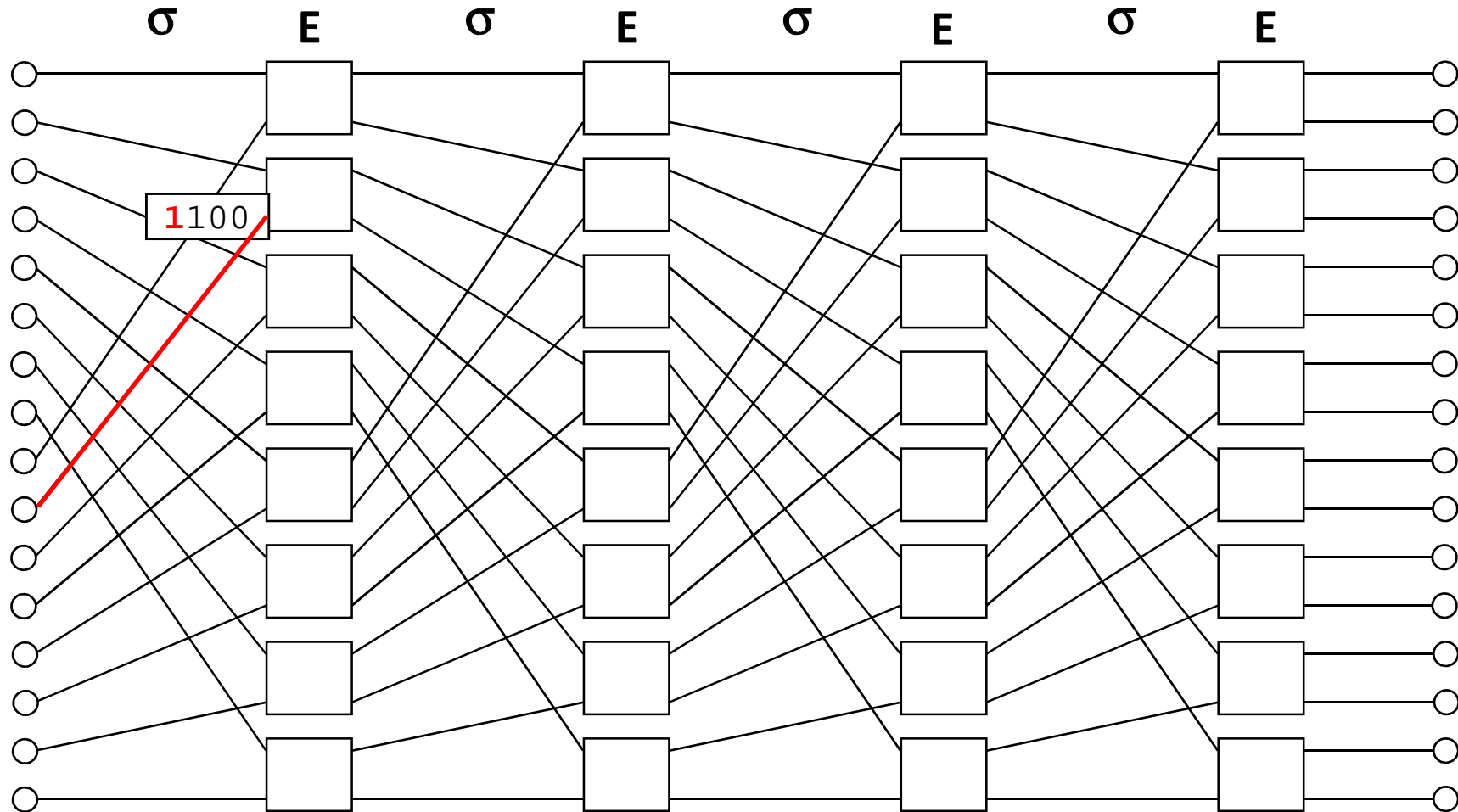
- Na základě adresy určení, tj binární reprezentace čísla výstupního portu (routing tag = destination port id) je možné řídit nastavení přepínačů v jednotlivých stupních sítě následovně. Přepínač v k-tém stupni sítě čte k-tý bit směrovacího návěští (počínaje od MSB směrem k LSB) a pokud se tento rovná:
  - 0 pak jdi na horní výstup
  - 1 pak jdi na dolní výstup



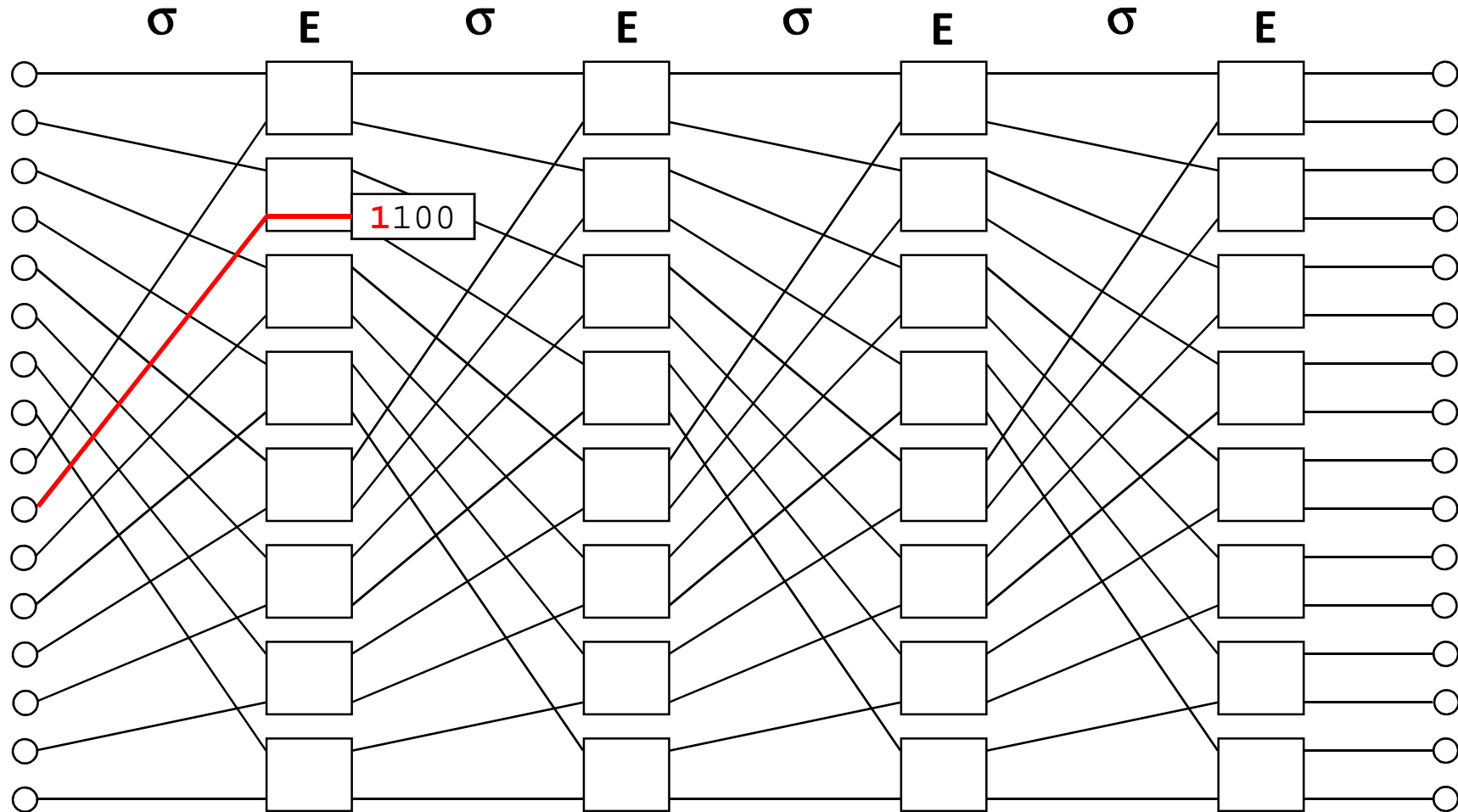
# Omega síť



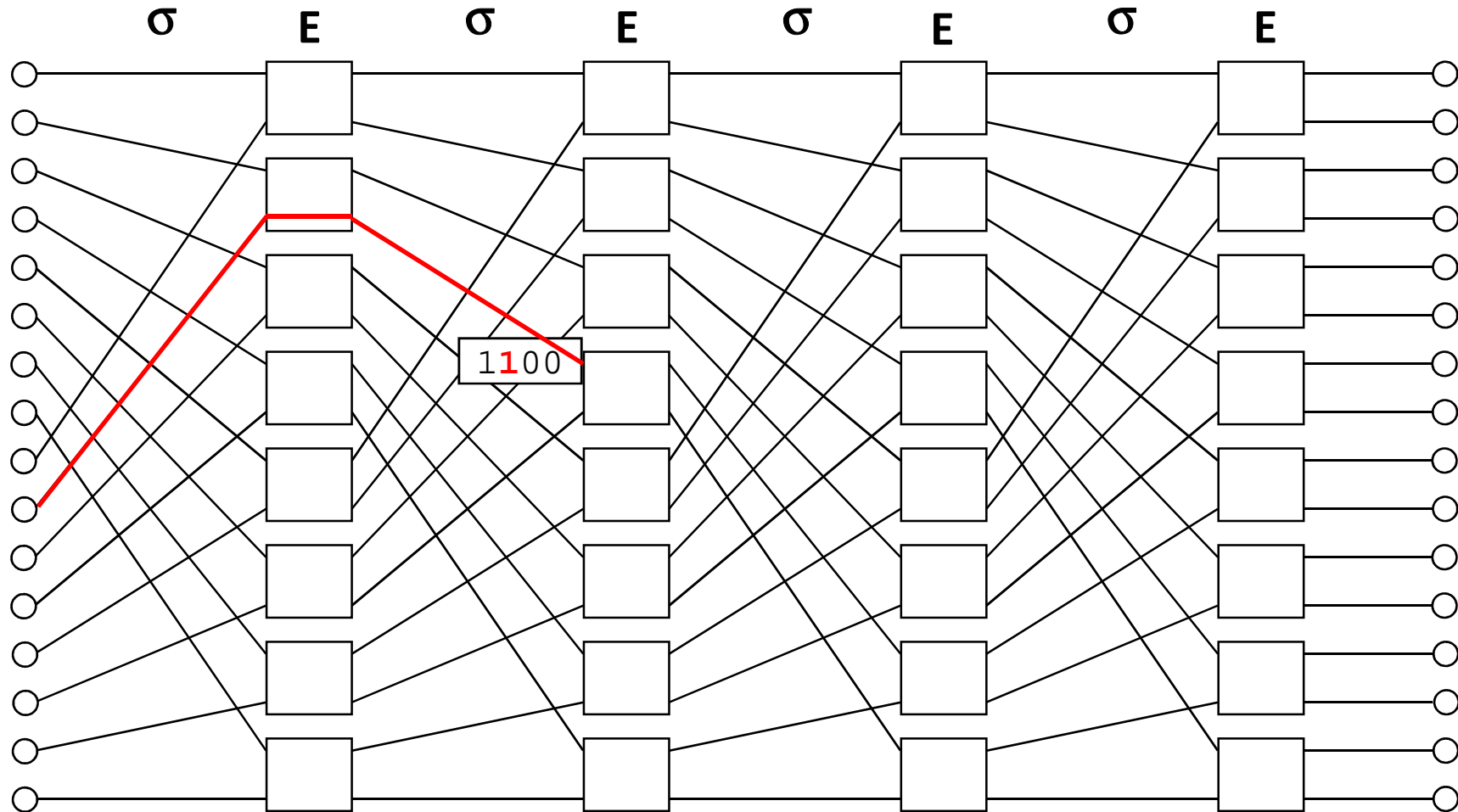
# Omega síť



# Omega síť

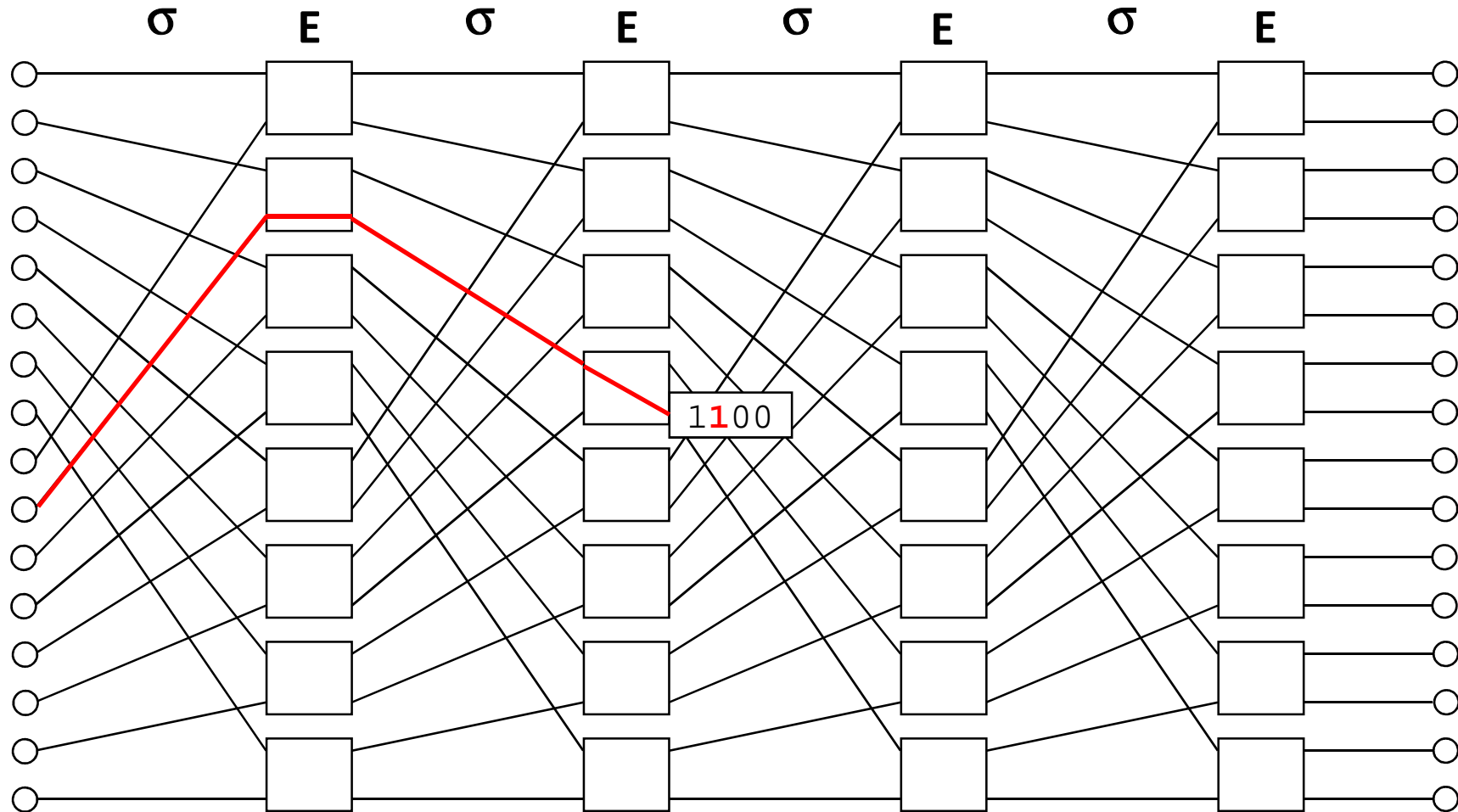


# Omega síť

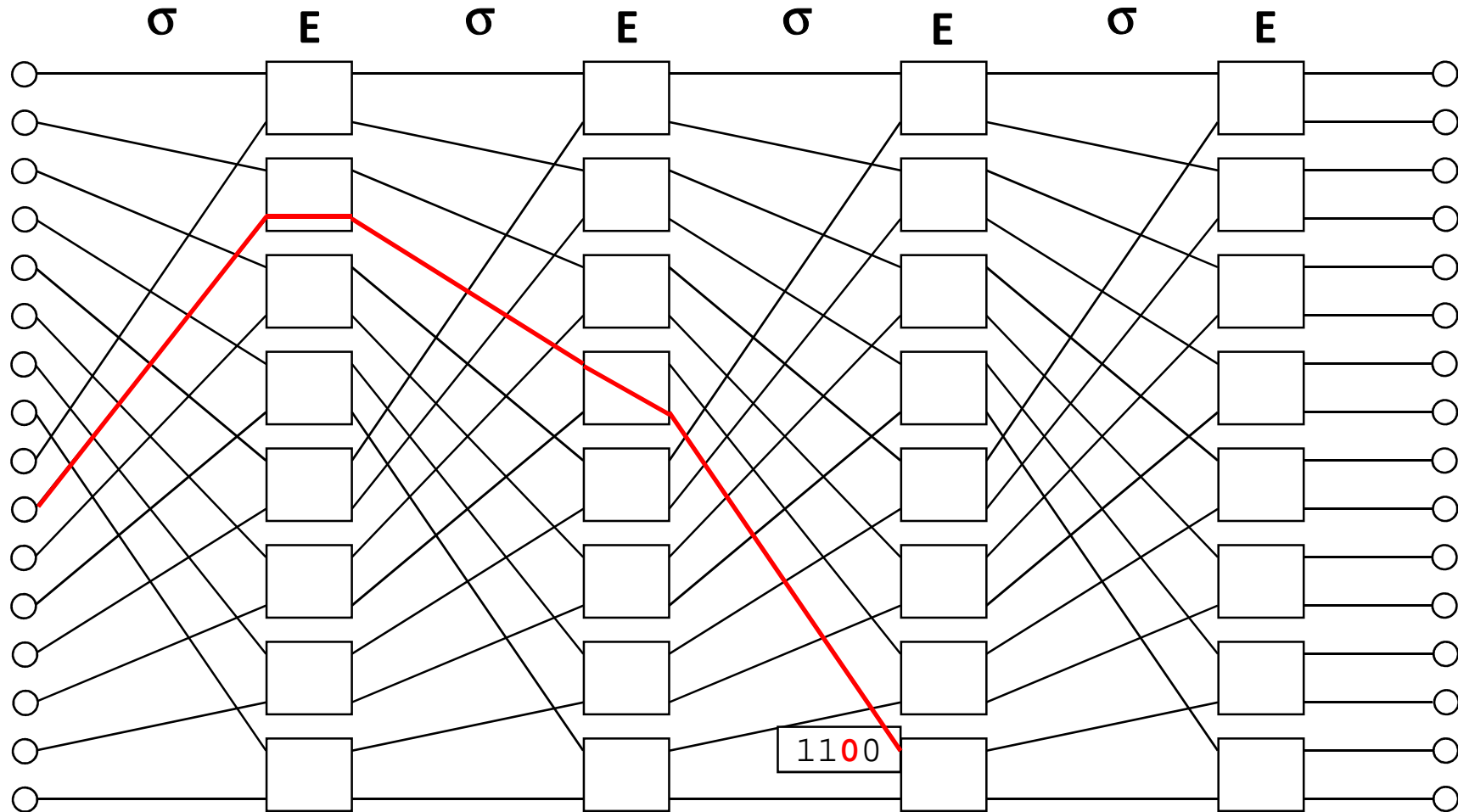




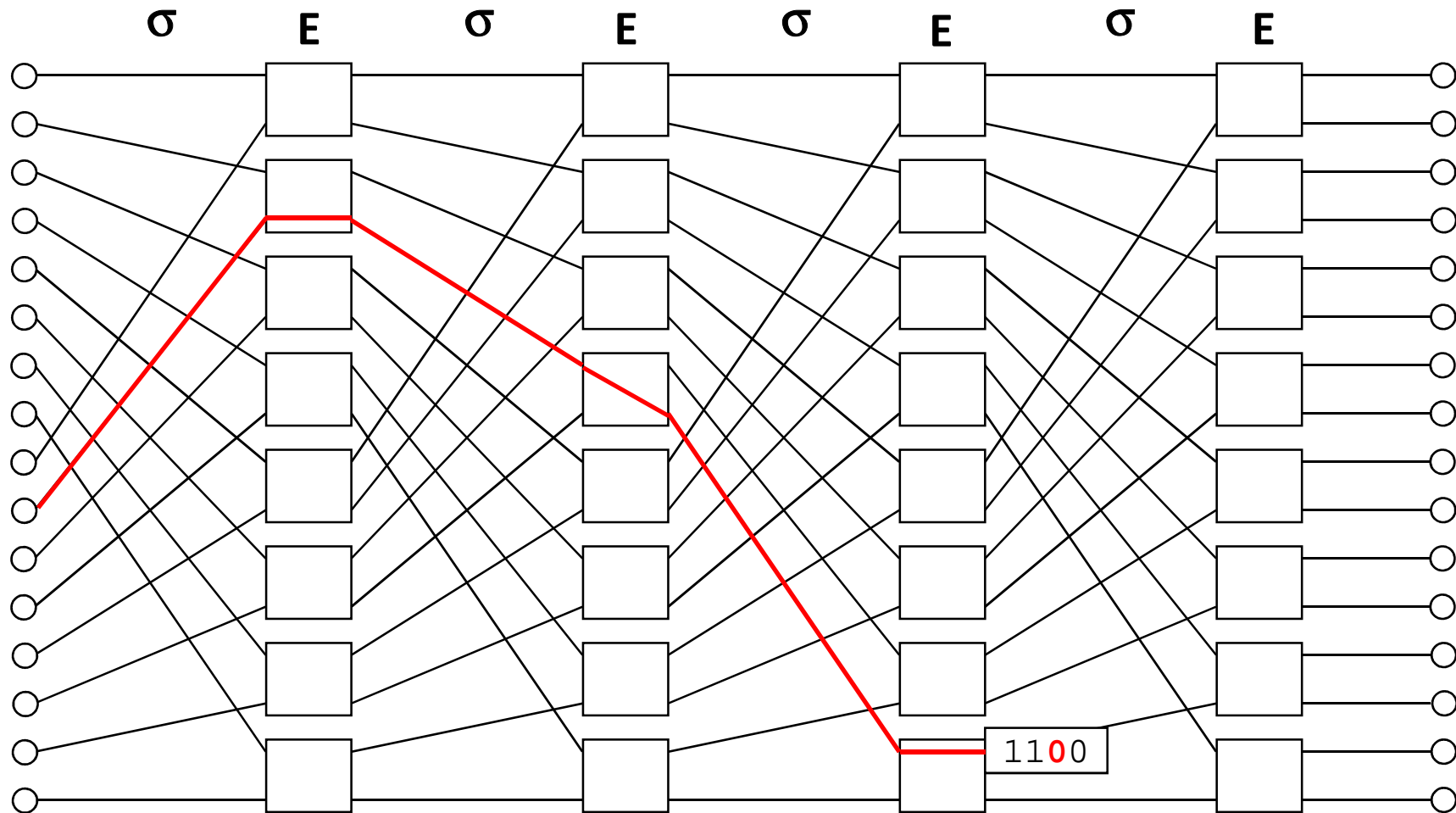
# Omega síť



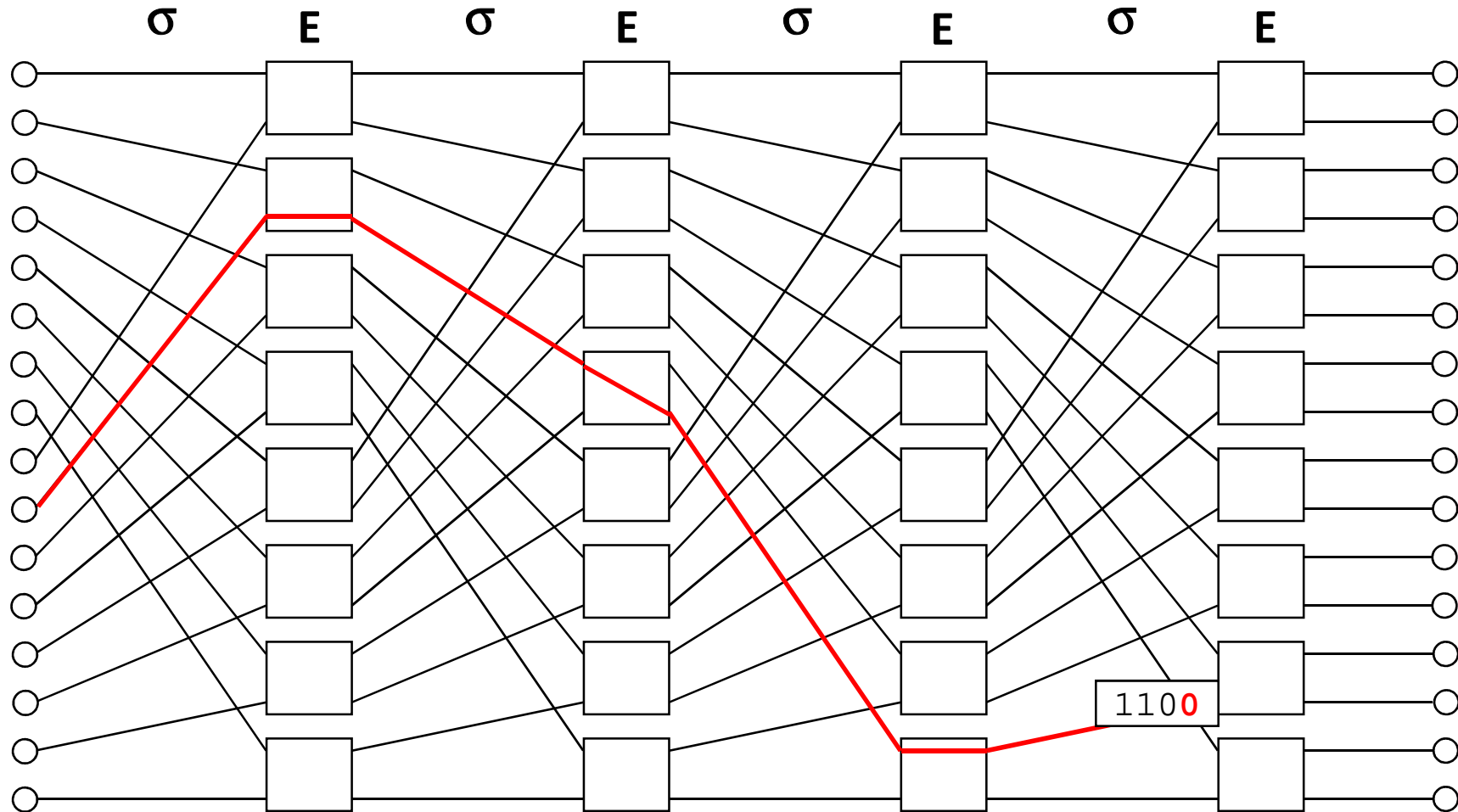
# Omega síť



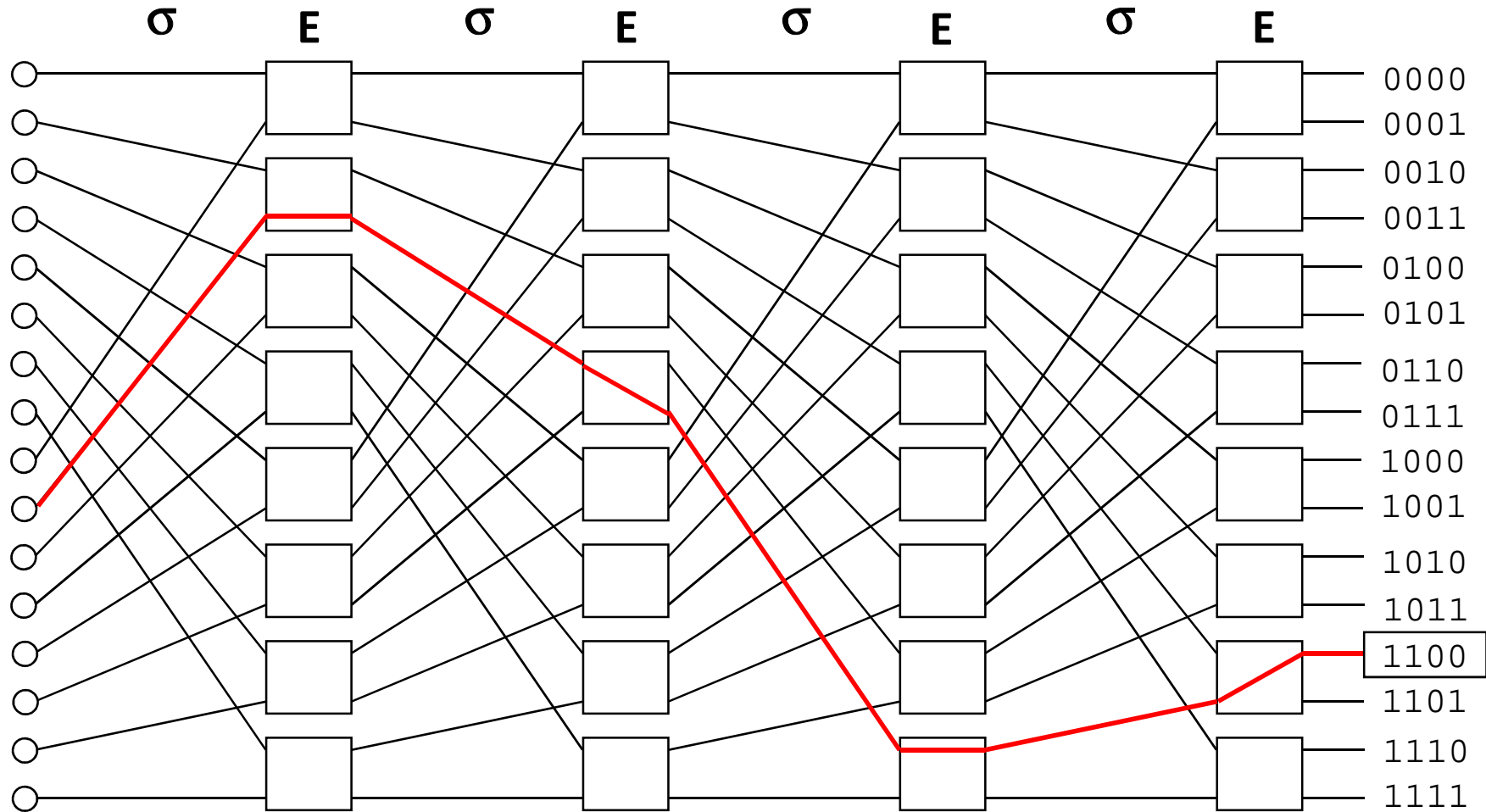
# Omega síť



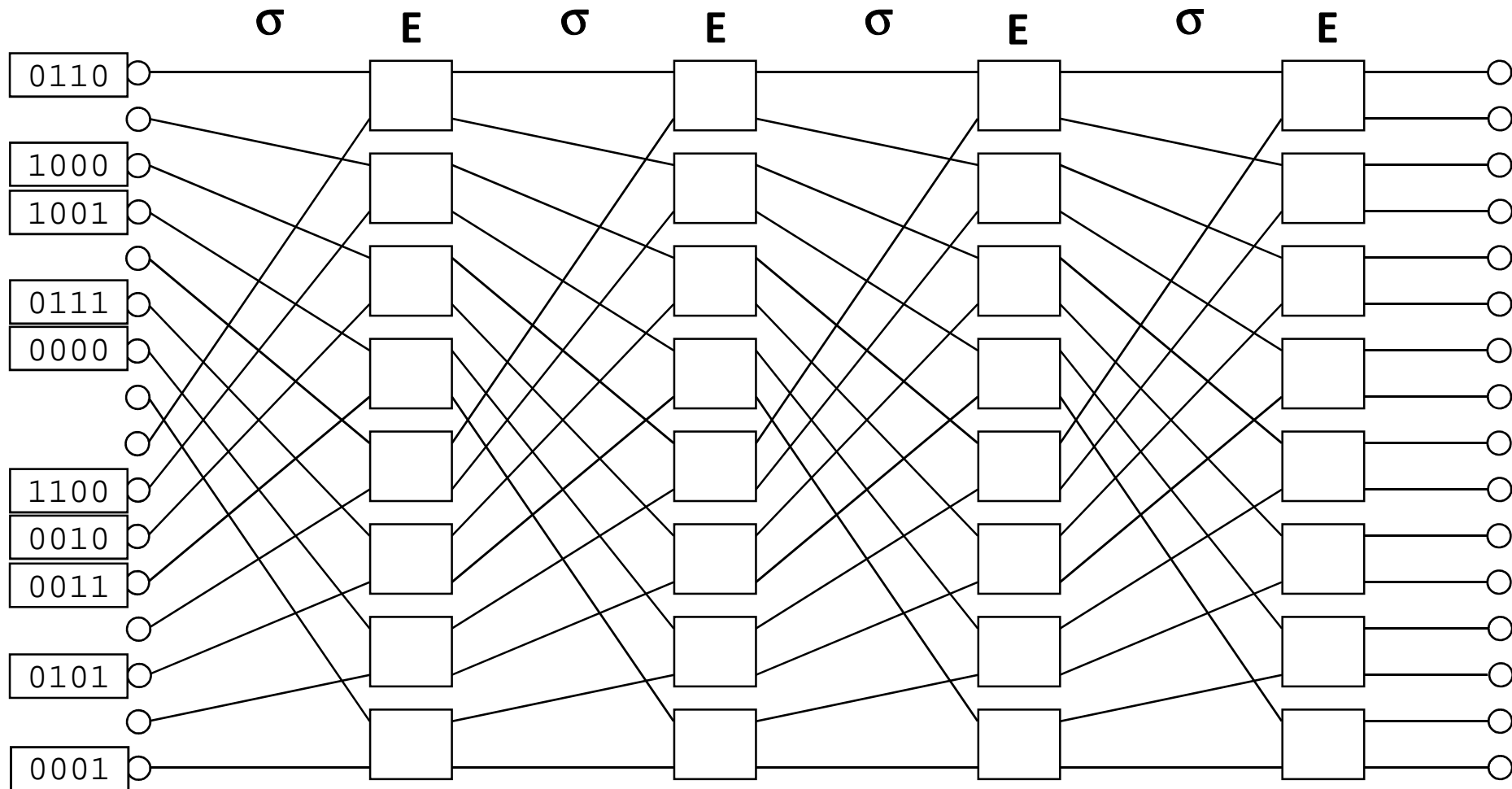
# Omega síť



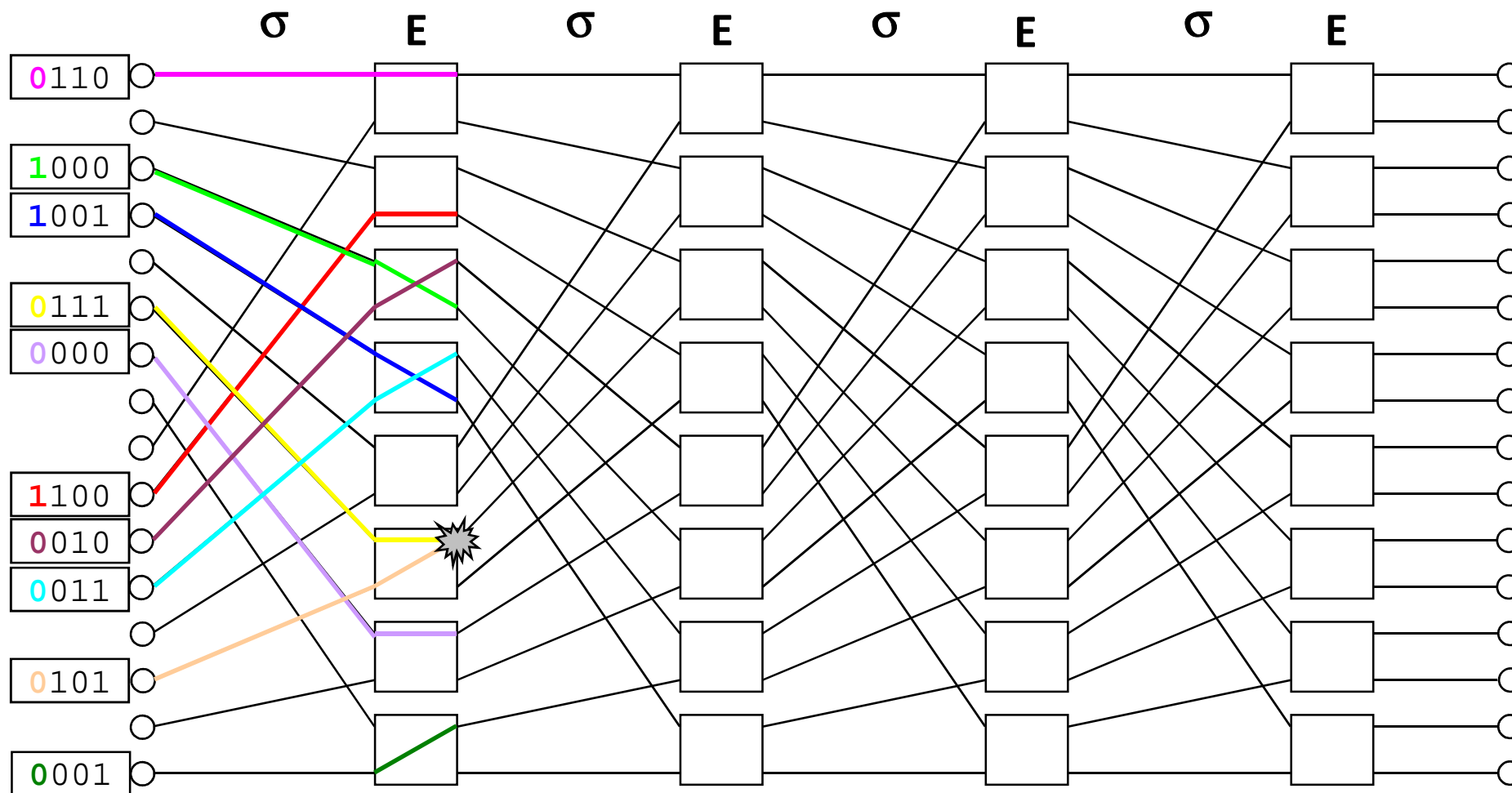
# Omega síť



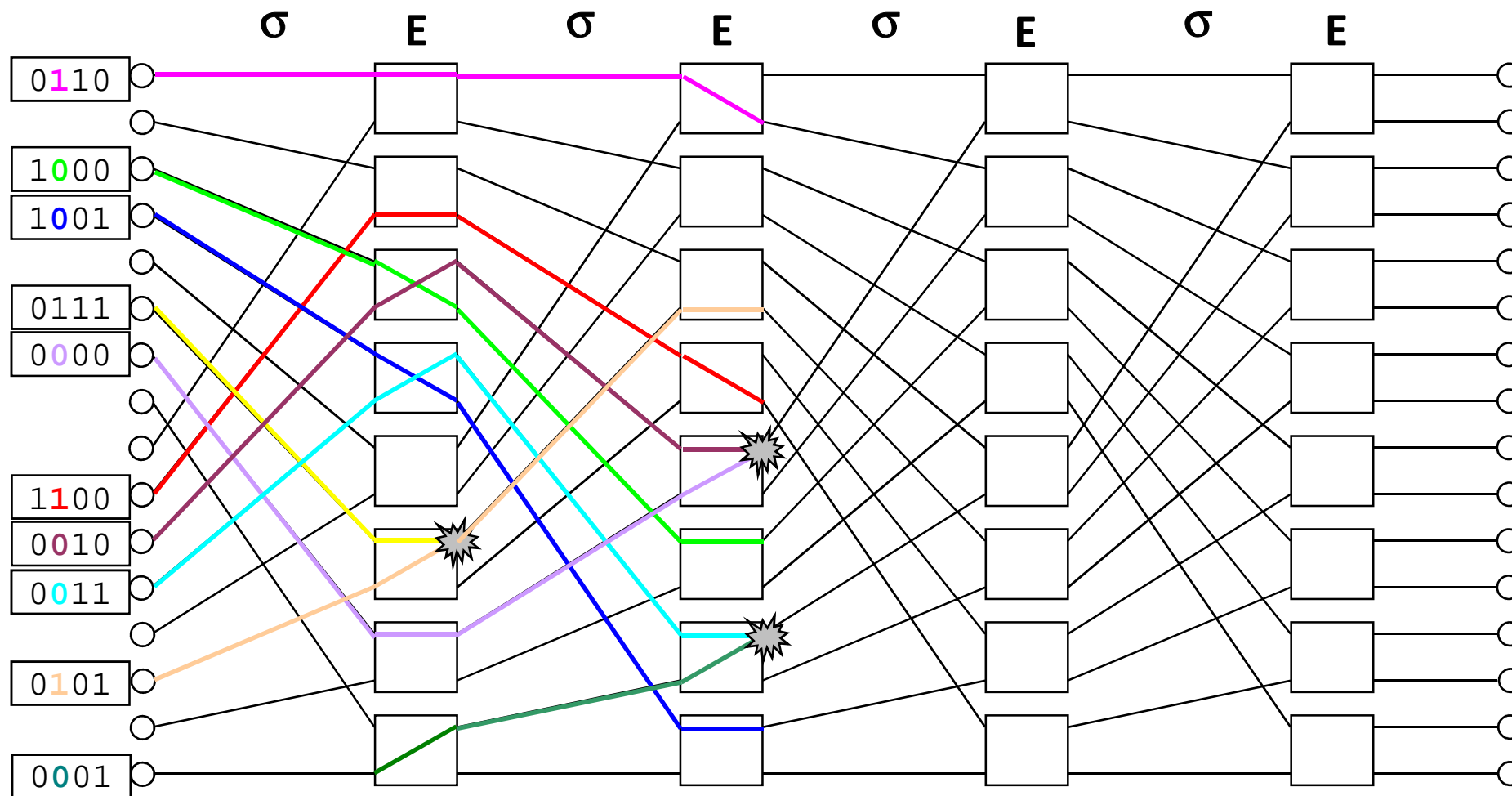
# Omega síť



# Omega síť



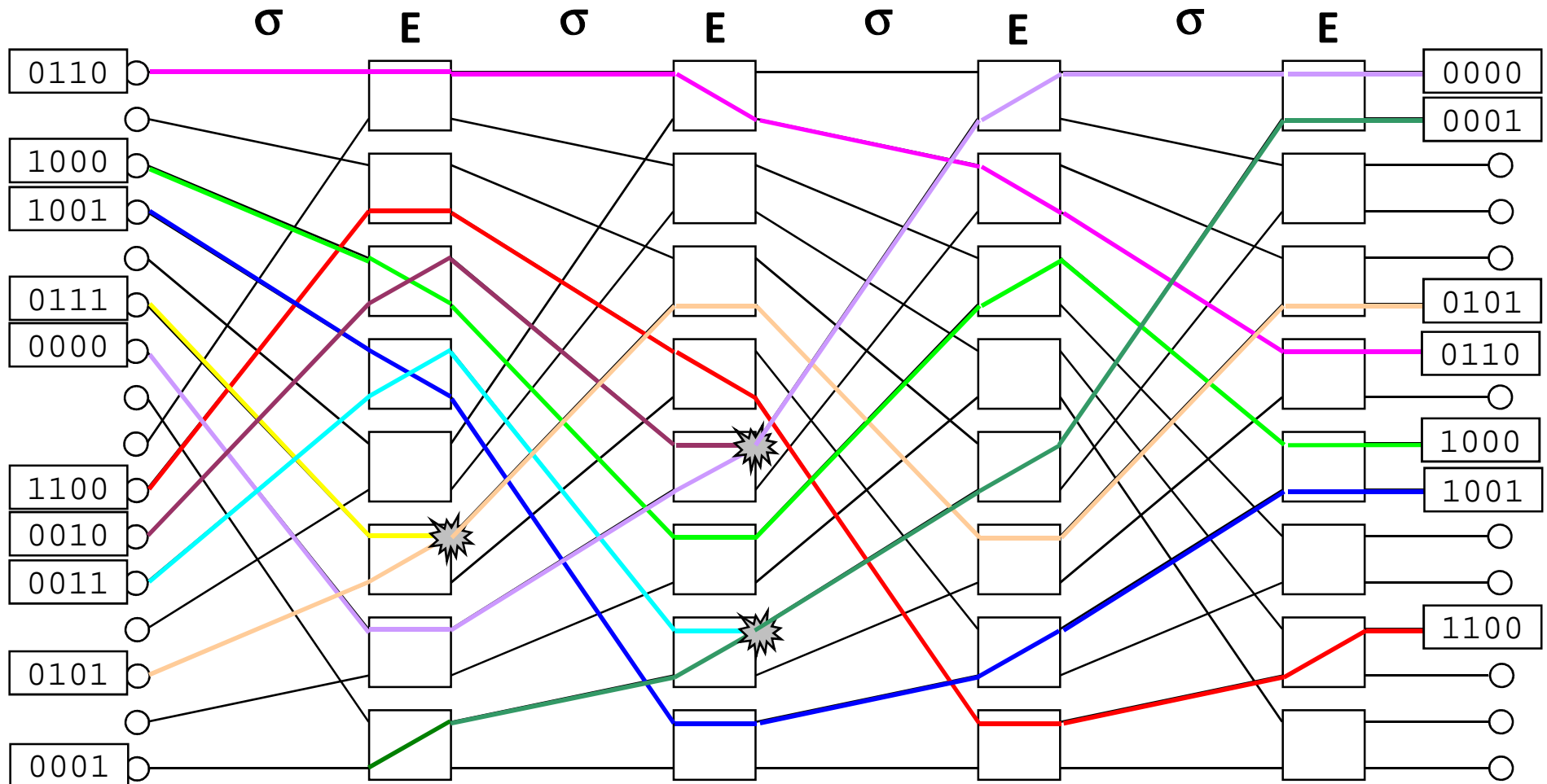
# Omega síť





# Omega síť

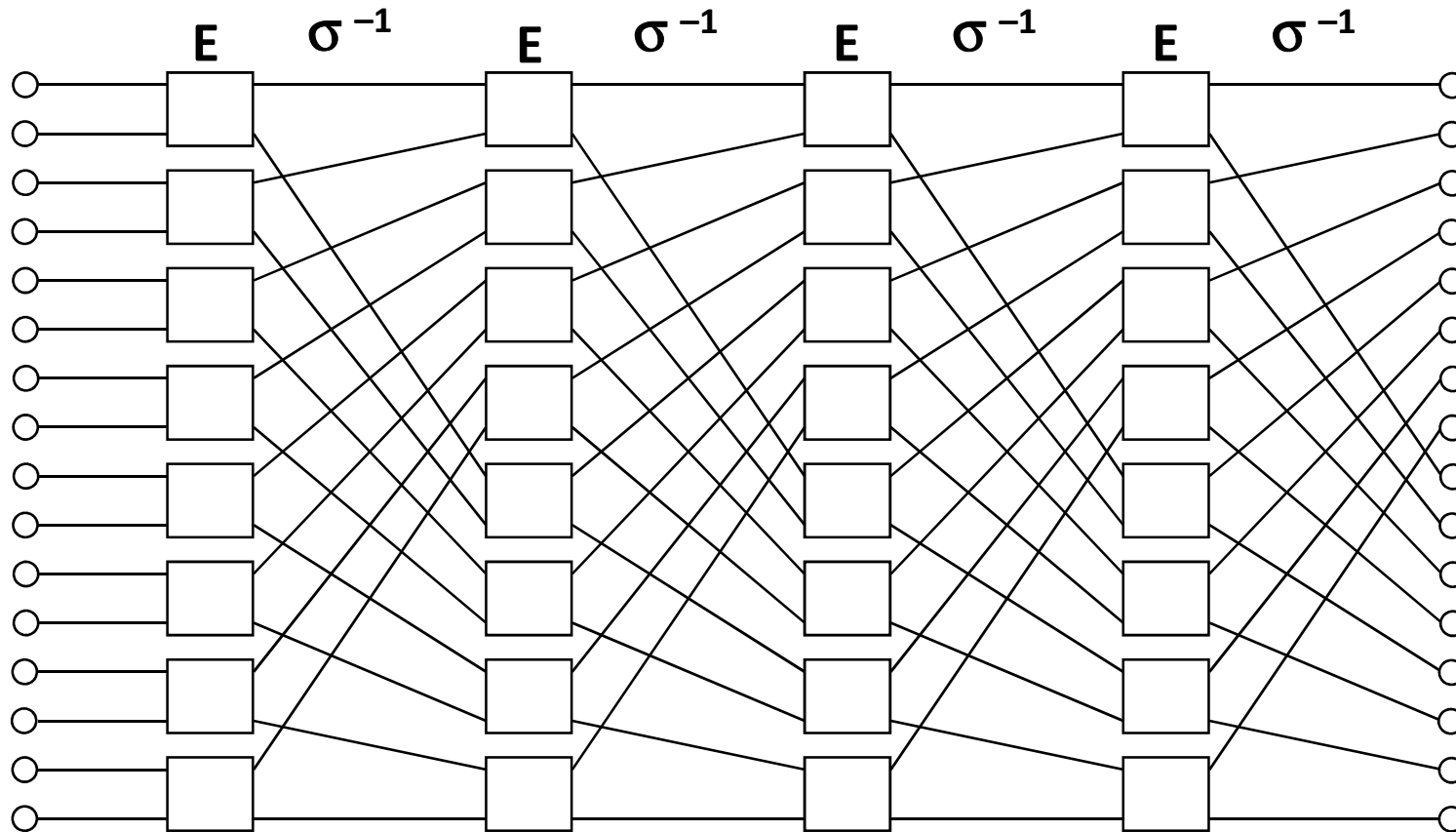
k výstupu dorazilo 7 z 10



## Inverzní omega síť

$$\Omega_N^{-1} = (\Omega_N)^{-1} = ((\sigma E)^n)^{-1} = ((\sigma E)^{-1})^n = (E^{-1} \sigma^{-1})^n = (E \sigma^{-1})^n$$

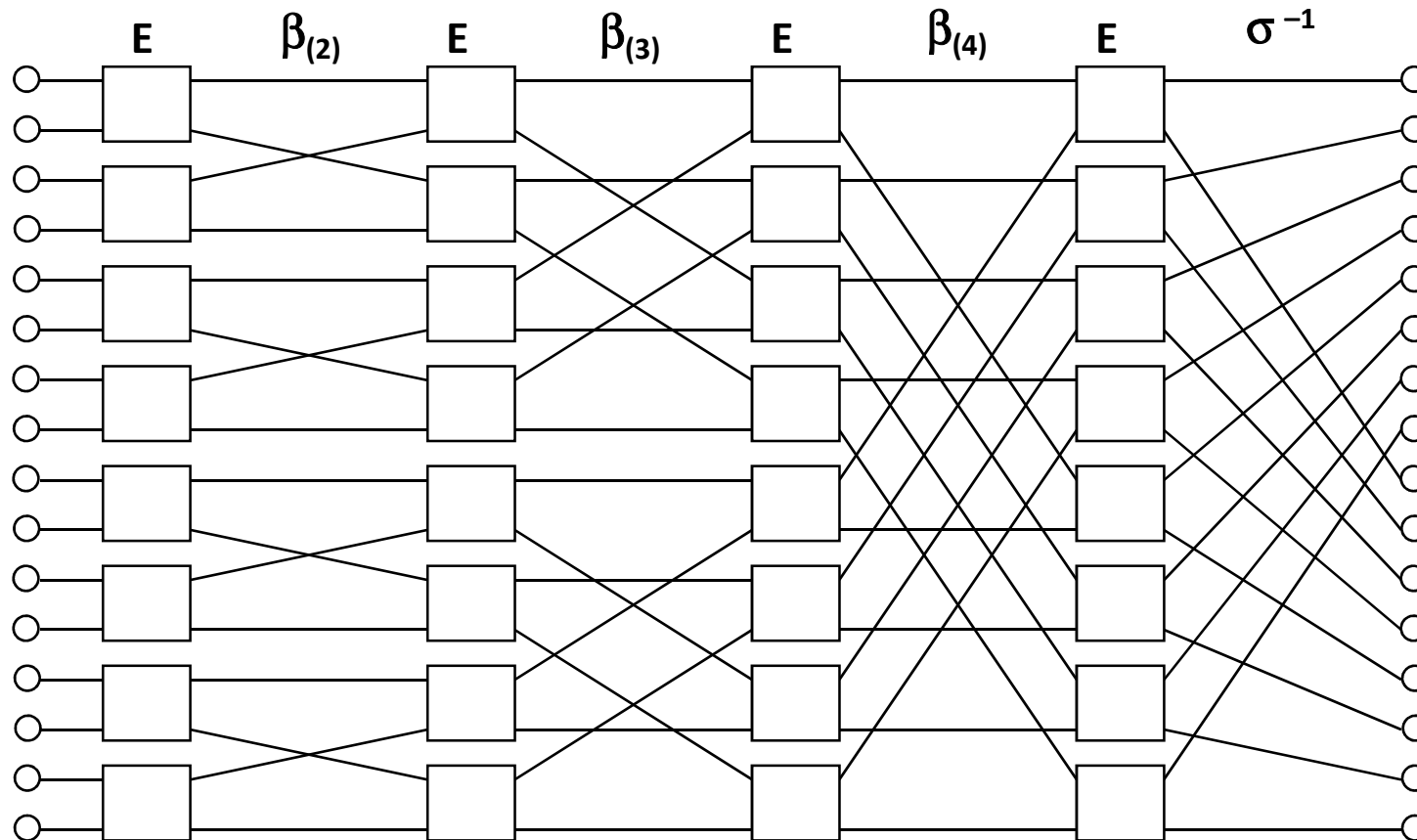
Takže  $\Omega_N^{-1} = (E \sigma^{-1})^n$



## Nepřímá binární n-kubická síť (Indirect binary cube)

$$\mathbf{C}_N = \mathbf{E}\beta_{(2)}\mathbf{E}\beta_{(3)}\mathbf{E}\beta_{(4)} \dots \mathbf{E}\beta_{(n)}\mathbf{E}\sigma^{-1}$$

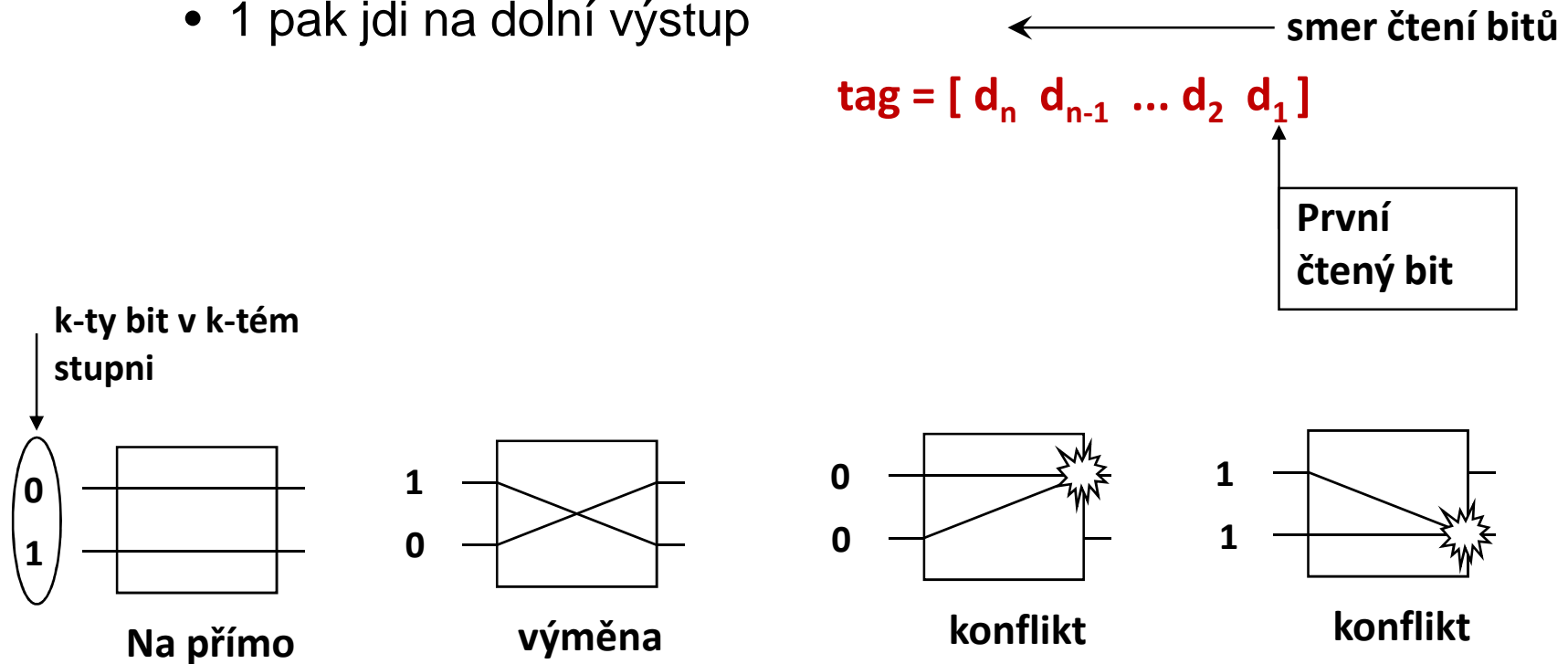
Napr. pro  $N=16$ ,  $n=4 \rightarrow \mathbf{C}_{16} = \mathbf{E}\beta_{(2)}\mathbf{E}\beta_{(3)}\mathbf{E}\beta_{(4)}\mathbf{E}\sigma^{-1}$



# Nepřímá binární n-kubická síť (Indirect binary cube)

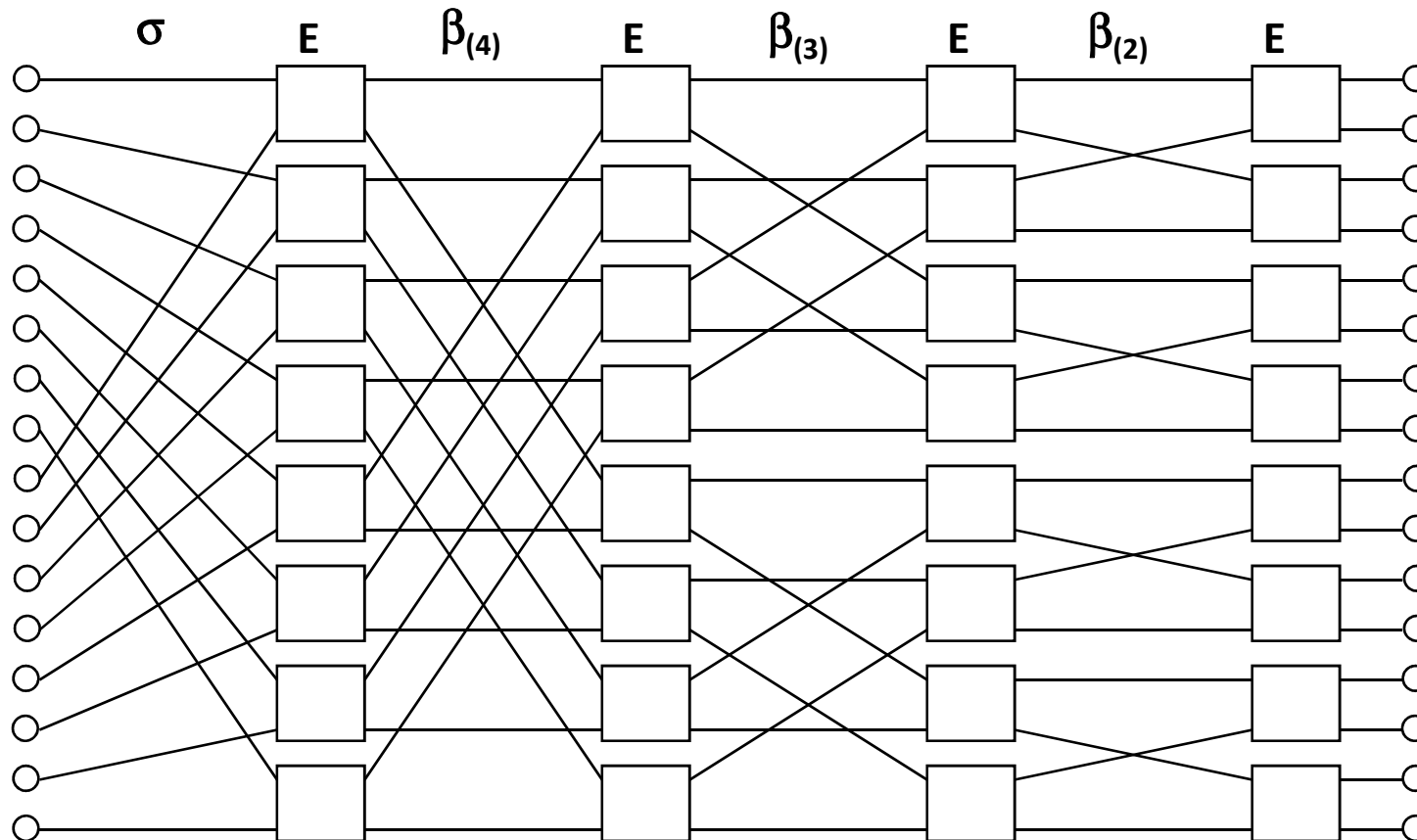
## Samosměrovací algoritmus:

- Přepínač v k-tém stupni sítě čte k-tý bit směrovacího návěští (počínaje LSB směrem k MSB) a pokud se tento rovná:
  - 0 pak jde na horní výstup
  - 1 pak jde na dolní výstup



## Inverzní nepřímá binární n-kubická síť

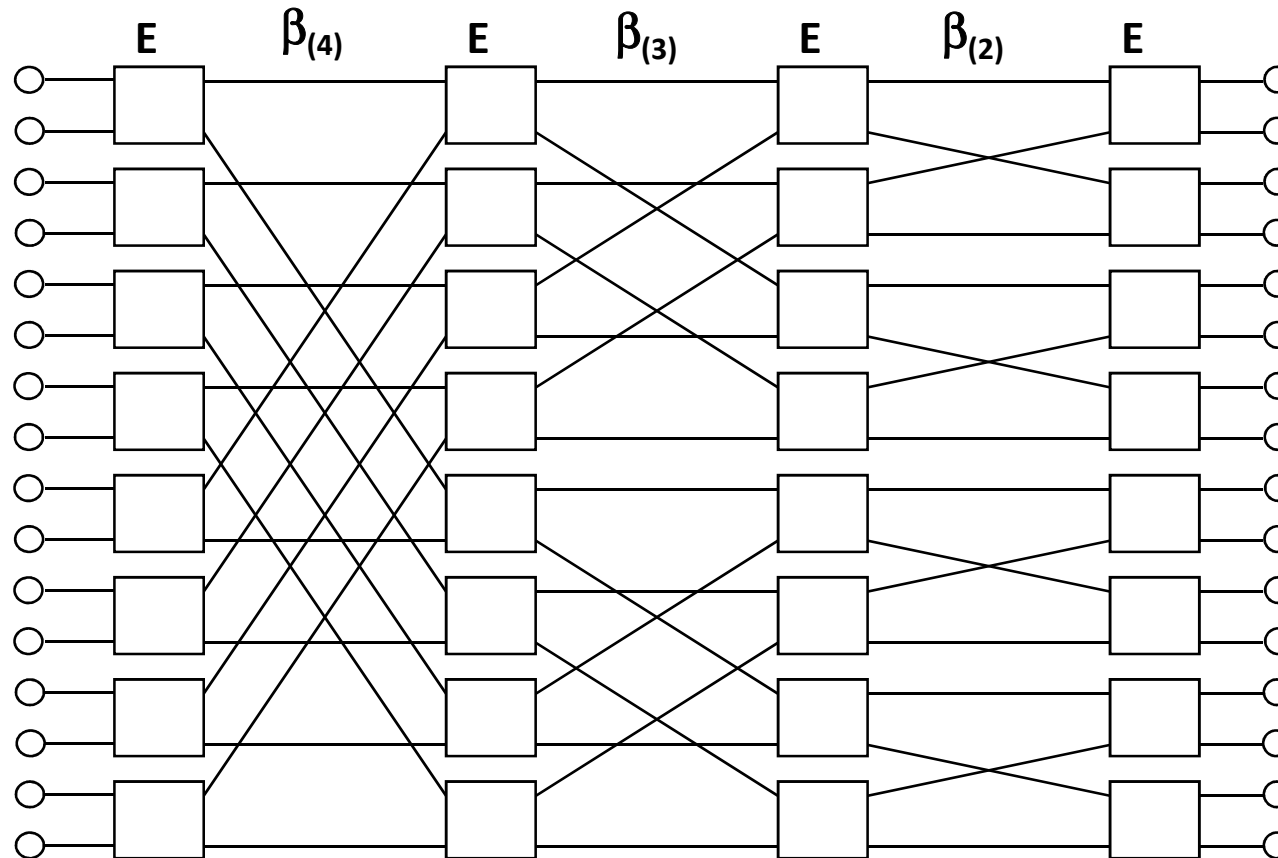
$$\mathbf{C}_N^{-1} = (\mathbf{E}\beta_{(2)}\mathbf{E}\beta_{(3)}\mathbf{E}\beta_{(4)} \dots \mathbf{E}\beta_{(n)}\mathbf{E}\sigma^{-1})^{-1} = \sigma \mathbf{E}\beta_{(n)} \dots \mathbf{E}\beta_{(4)} \mathbf{E}\beta_{(3)} \mathbf{E}\beta_{(2)} \mathbf{E}$$



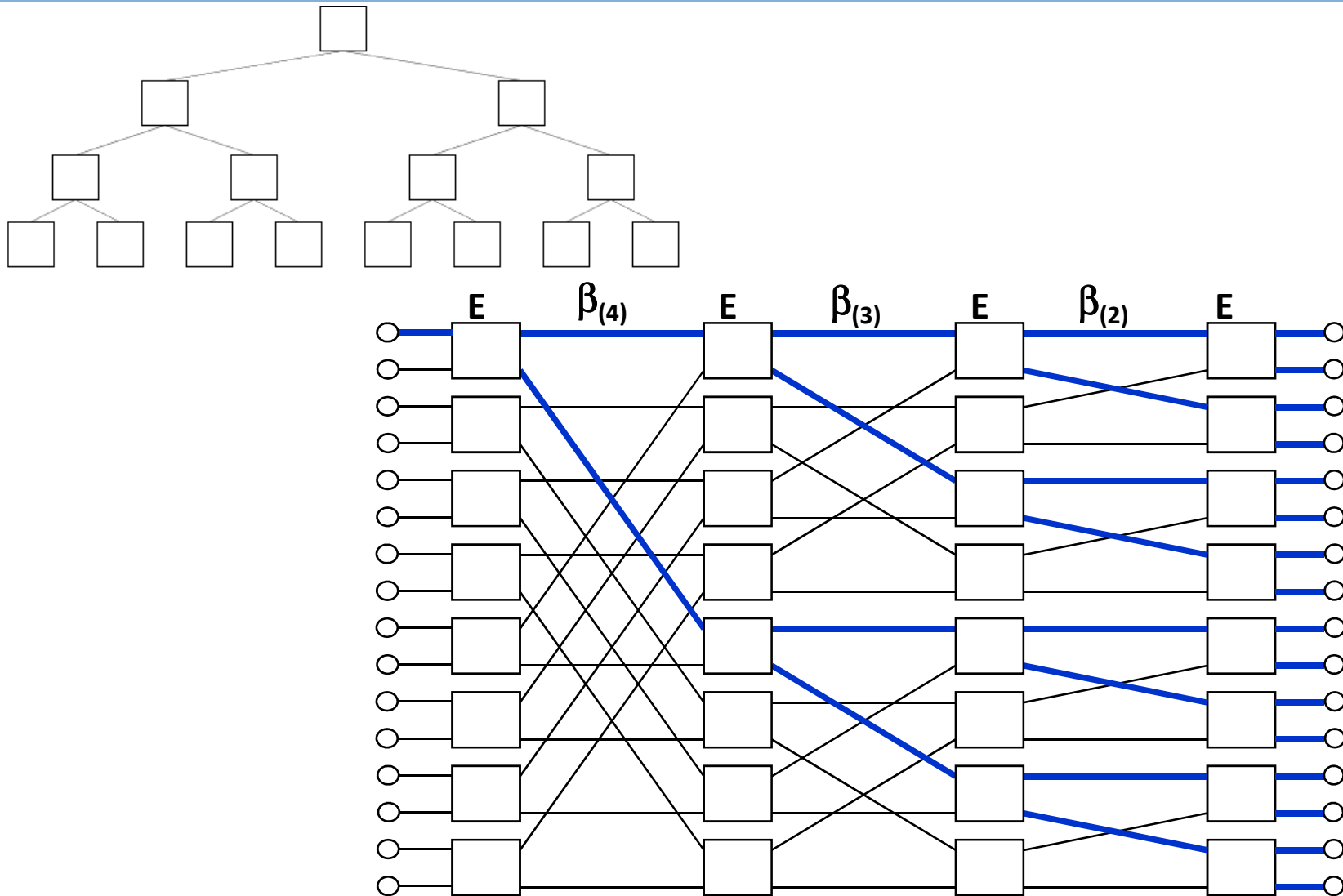
# Butterfly network

$$C_N^{-1} = (E\beta_{(2)}E\beta_{(3)}E\beta_{(4)} \dots E\beta_{(n)}E\sigma^{-1})^{-1} = \sigma E\beta_{(n)} \dots E\beta_{(4)} E\beta_{(3)} E\beta_{(2)} E$$

$$B = E\beta_{(n)} \dots E\beta_{(4)} E\beta_{(3)} E\beta_{(2)} E$$



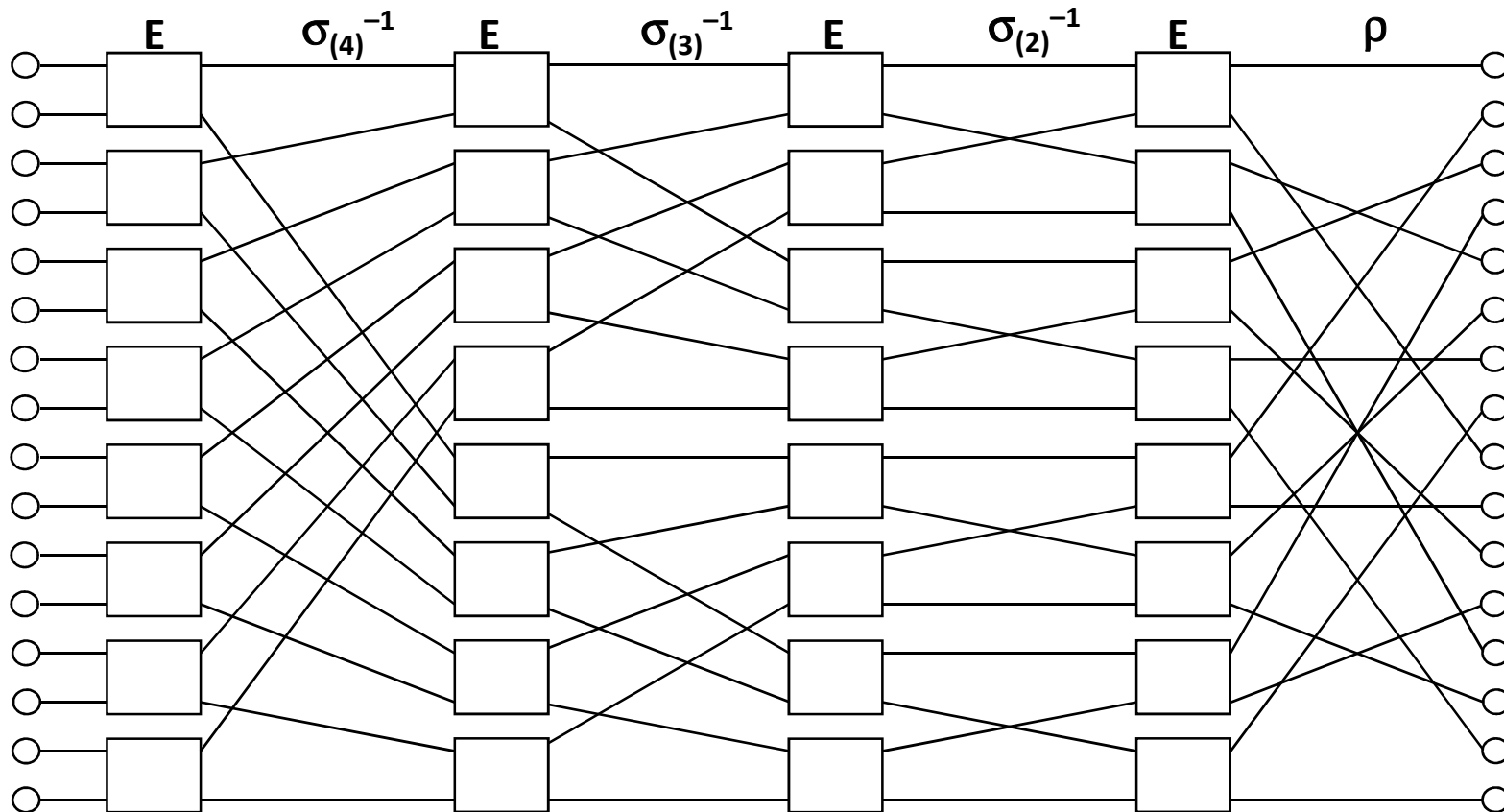
# Butterfly network



## R-sít' (R-network)

$$R_N = E\sigma_{(n)}^{-1} E\sigma_{(n-1)}^{-1} \dots E\sigma_{(2)}^{-1} E\rho$$

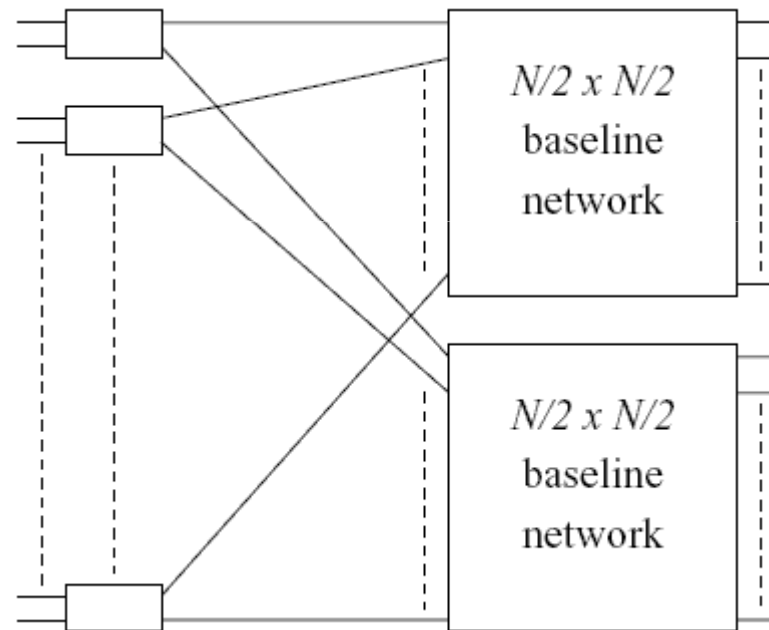
Napr. pro  $N=16$ ,  $n=4 \rightarrow R_{16} = E\sigma_{(4)}^{-1} E\sigma_{(3)}^{-1} E\sigma_{(2)}^{-1} E\rho$





## Baseline síť (Baseline network)

- Topologie baseline sítě je definována rekurzivní způsobem podle následujícího obrázku:



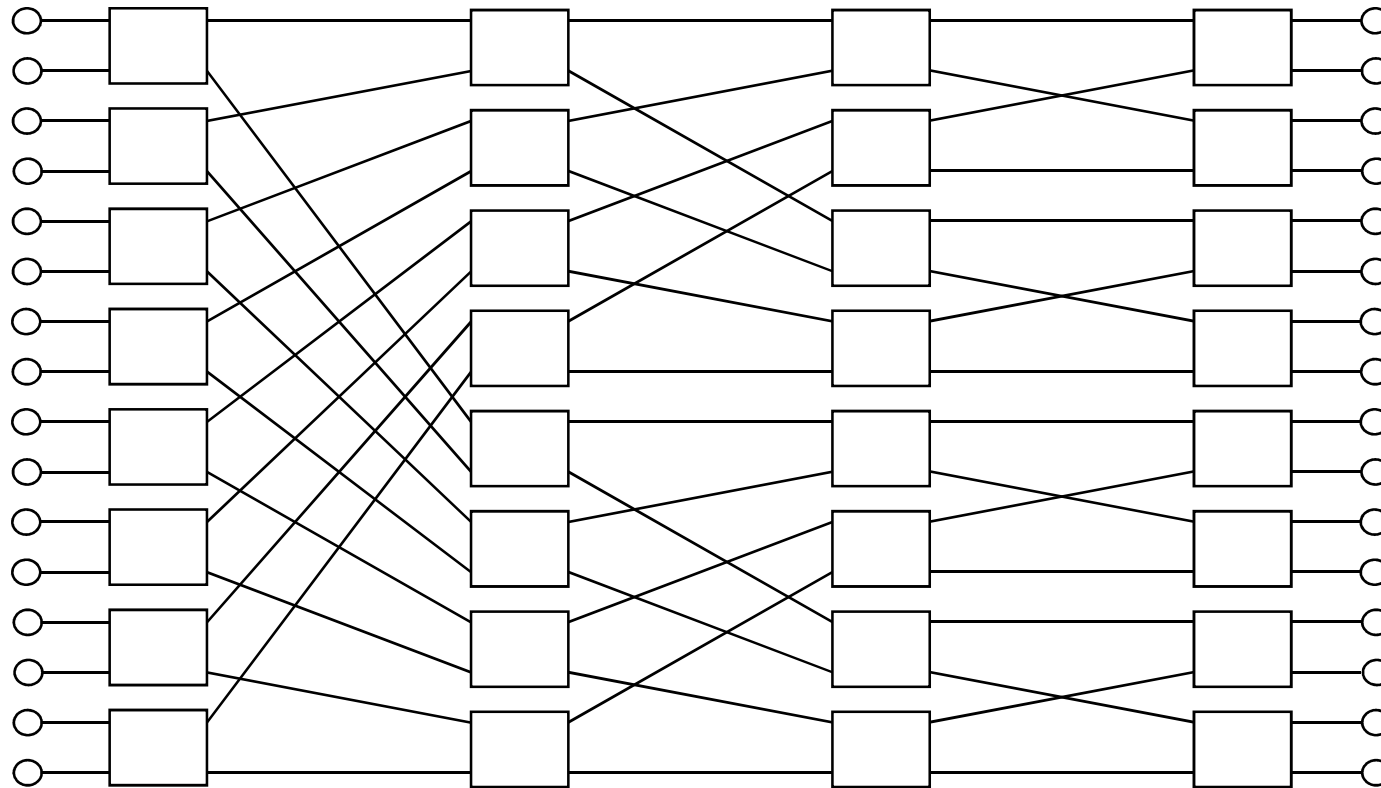
- Spojovací funkce mezi dvěma stupni je reverzní míchání ( $\sigma^{-1}$ ).

## Baseline síť (Baseline network)

- Rekurzivní proces se ukončí pokud dosáhneme elementy 2x2.  
Rozepsáním rekurze můžeme získat zápis baseline sítě:

$$\mathbf{Baseline}_N = \mathbf{E}\sigma_{(n)}^{-1} \mathbf{E}\sigma_{(n-1)}^{-1} \dots \mathbf{E}\sigma_{(2)}^{-1} \mathbf{E}, \quad \text{kde } n = \log_2 N$$

- Např. pro  $N = 16$ ,  $n = 4$ :

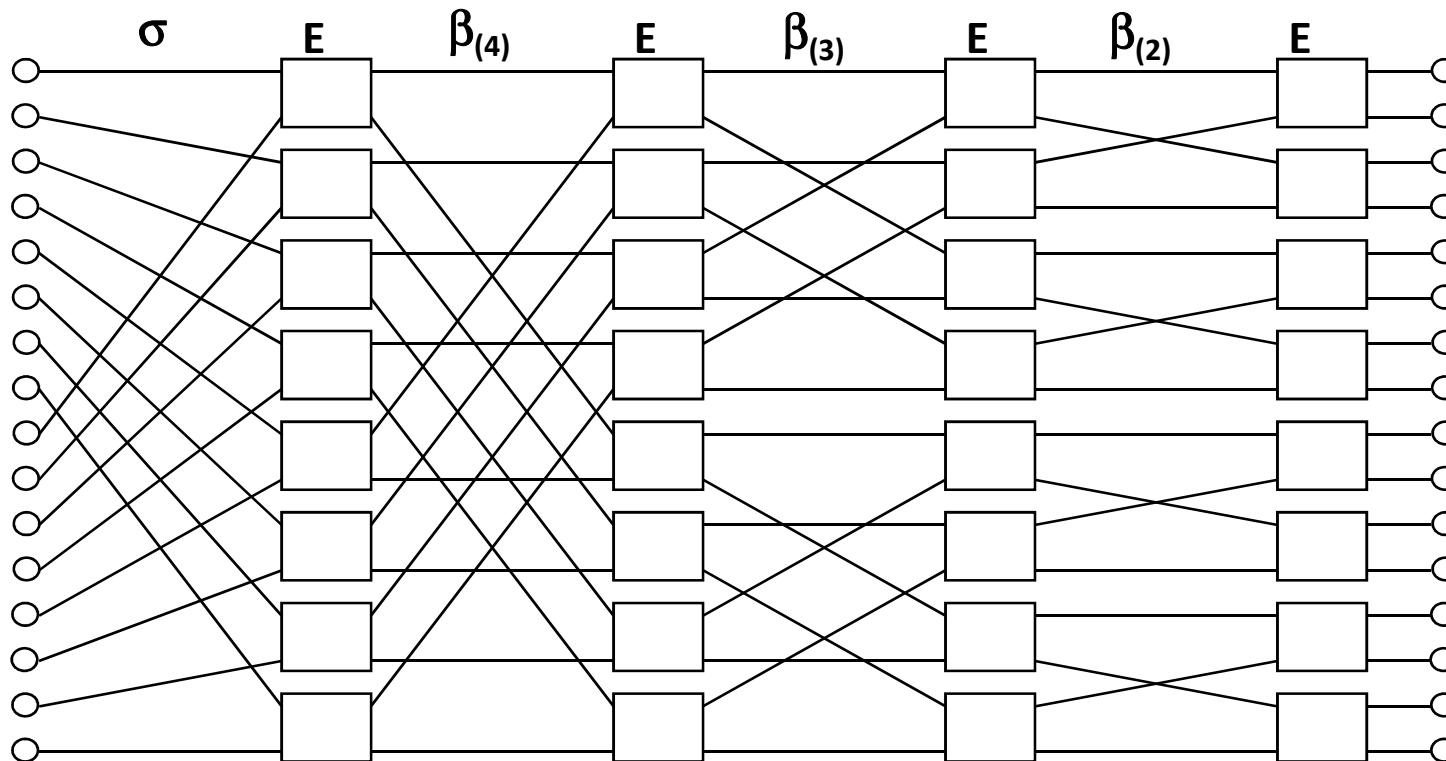


## Banyan network

- Banyan síť, někdy také nazývaná generalizovaná kubická síť (generalized cube network) je síť definována předpisem:

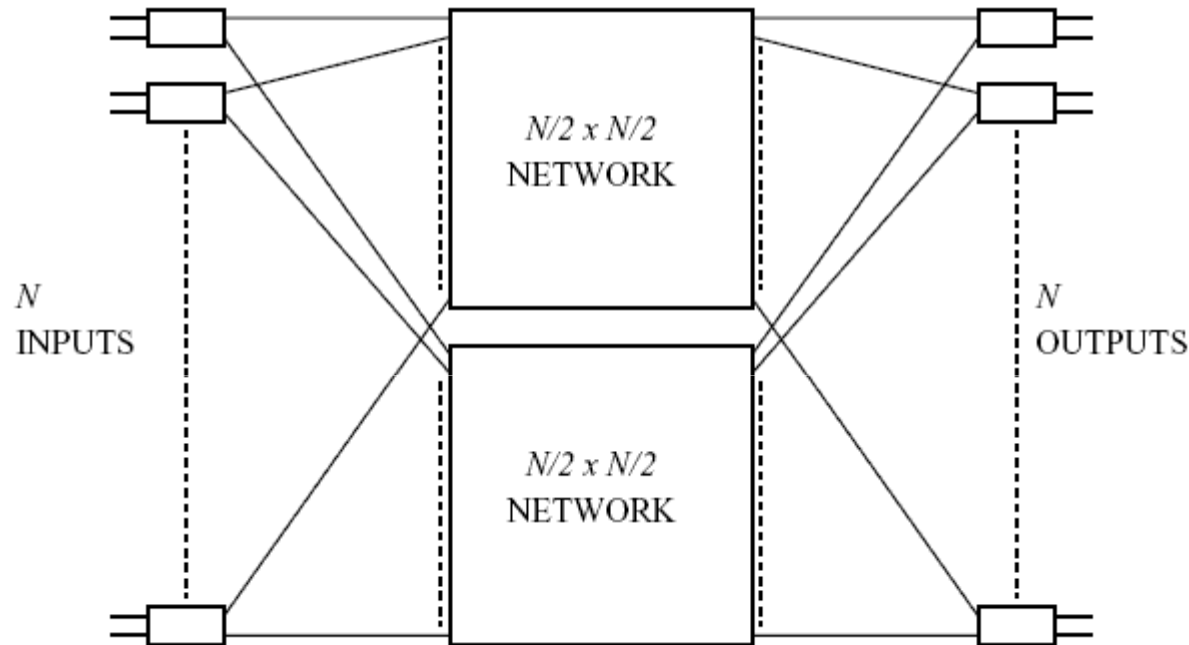
$$\mathbf{Banyan}_N = \sigma \mathbf{E} \beta_{(n)} \dots \mathbf{E} \beta_{(4)} \mathbf{E} \beta_{(3)} \mathbf{E} \beta_{(2)} \mathbf{E}, \quad \text{kde } n = \log_2 N$$

a tedy se jedná o inverzní nepřímou binární n-kubickou síť.



## Benešova síť (Beneš network)

- Rekurzivní definice Benešovy sítě:

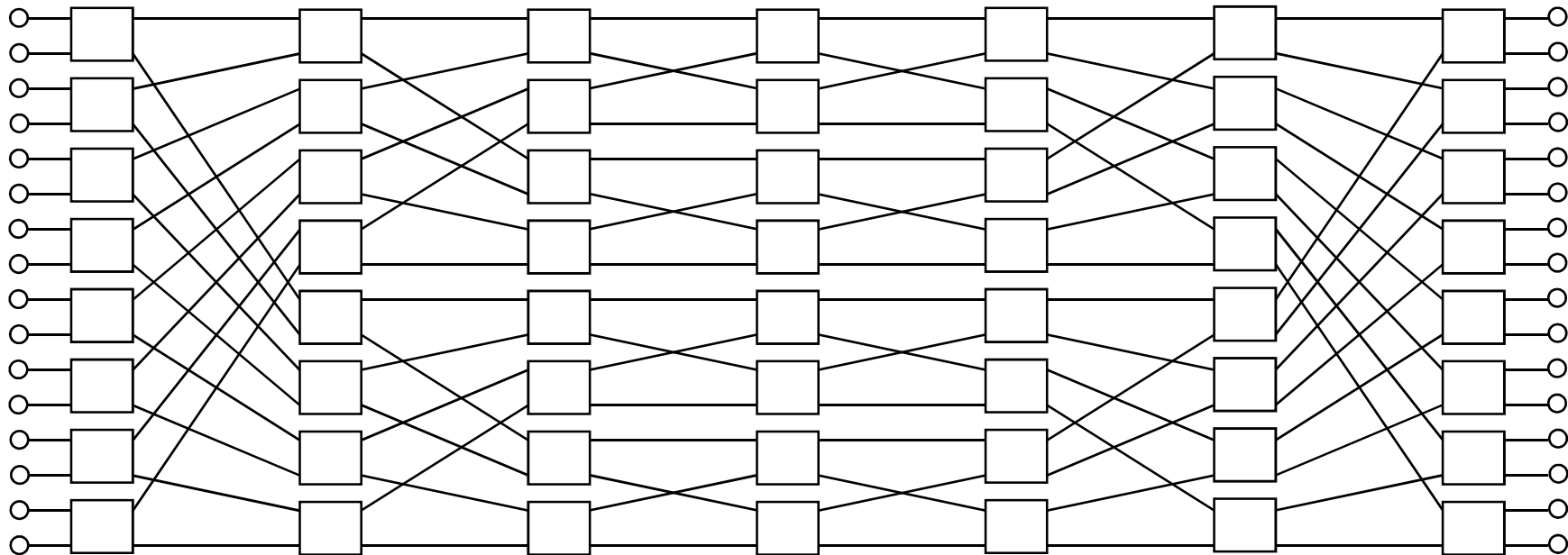


- Když je rekurzivní proces ukončen síť má  $2\log_2 N - 1$  stupňů a v každém stupni  $N/2$  spojovacích elementů. Benešova síť patří mezi **rekonfigurovatelné sítě bez blokování** a proto vyžaduje algoritmus pro rekonfiguračný proces - centrální řízení. Existují také samosměrovací algoritmy částečně potlačující blokování v síti.

## Benešova síť (Beneš network)

$$\mathbf{Beneš}_N = \mathbf{E}\sigma_{(k)}^{-1} \mathbf{E}\sigma_{(k-1)}^{-1} \mathbf{E} \dots \sigma_{(2)}^{-1} \mathbf{E} \sigma_{(2)} \dots \mathbf{E}\sigma_{(k-1)} \mathbf{E}\sigma_{(k)} \mathbf{E},$$

- Např. pro  $N = 16$ , stupňů:  $n = 2\log_2 N - 1 = 7$ :



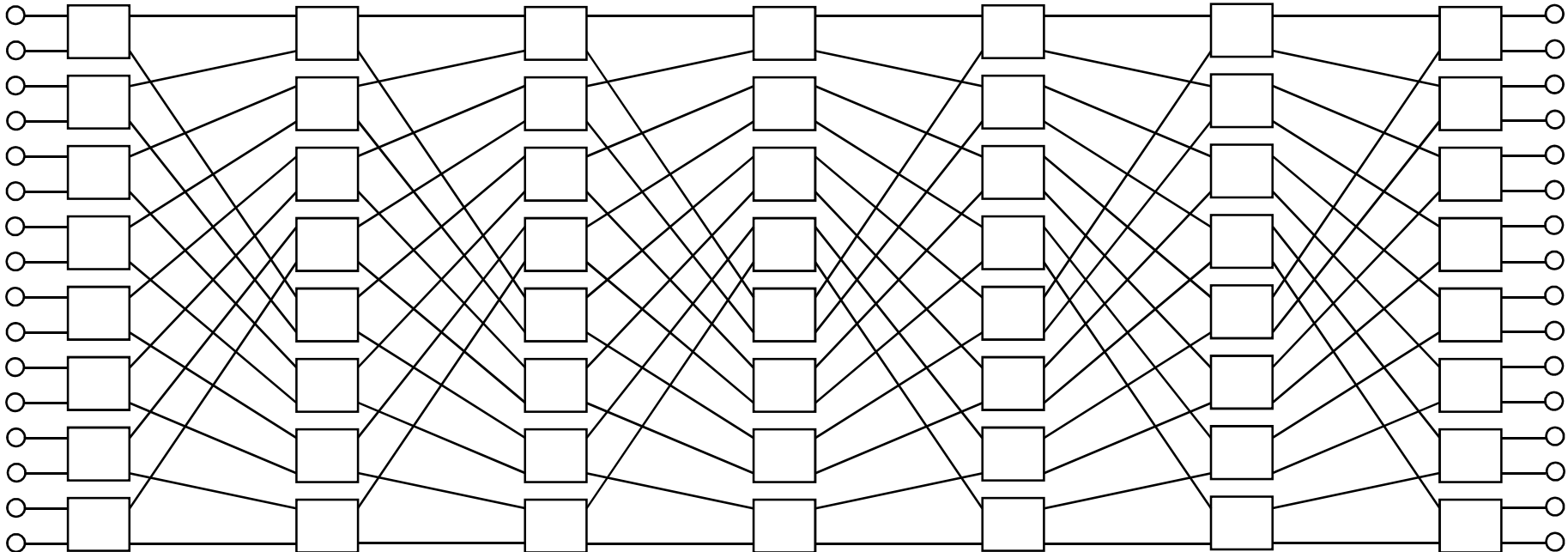
- Na Benešovu síť se také můžeme dívat jako na zařazení dvou baseline sítí za sebou (baseline síť a reverzní baseline síť). Proto bývají Benešově sítě někdy také nazývané "sériové baseline sítě".

## Benešova síť (Beneš network)

Alternativní definice:

$$\text{Beneš}_N = E\sigma^{-1} E\sigma^{-1} E \dots \sigma^{-1} E\sigma \dots E\sigma E\sigma E = (E\sigma^{-1})^{k-1} E (\sigma E)^{k-1},$$

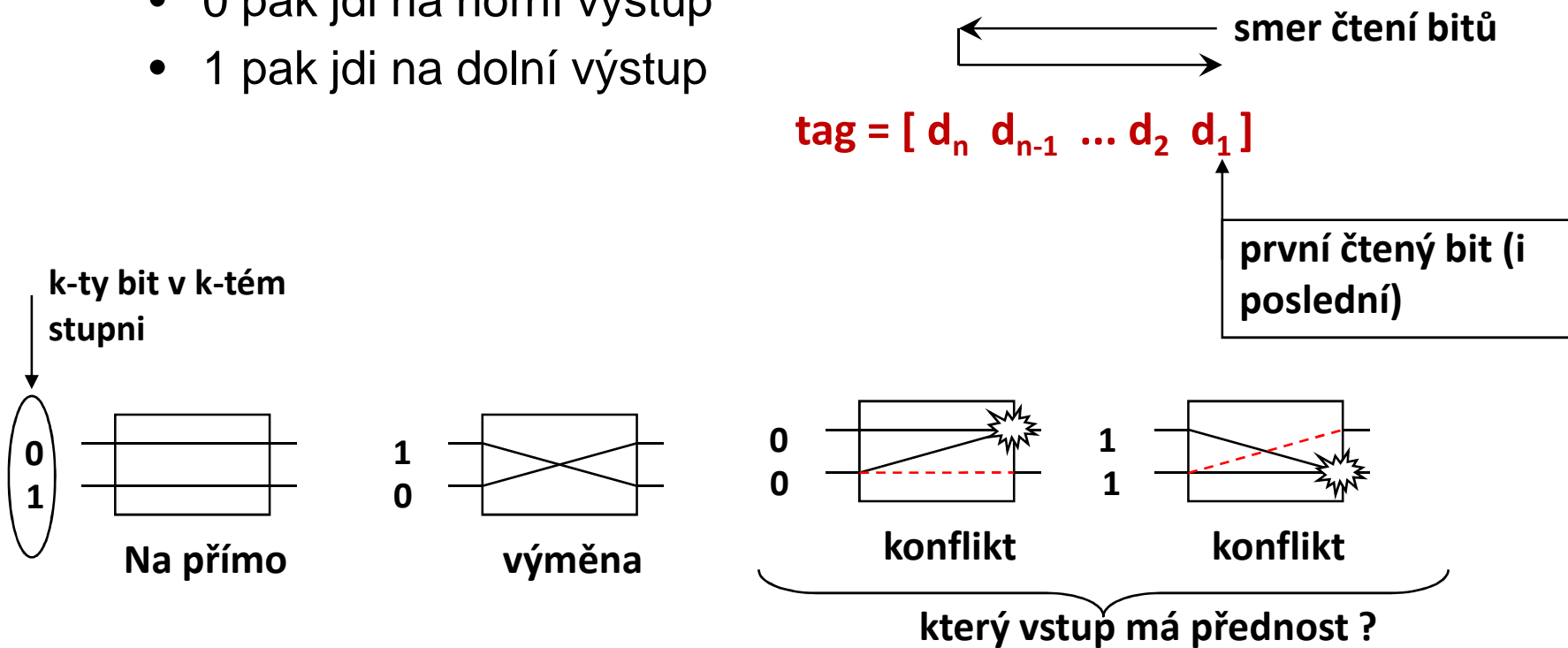
- kde  $k = \log_2 N$



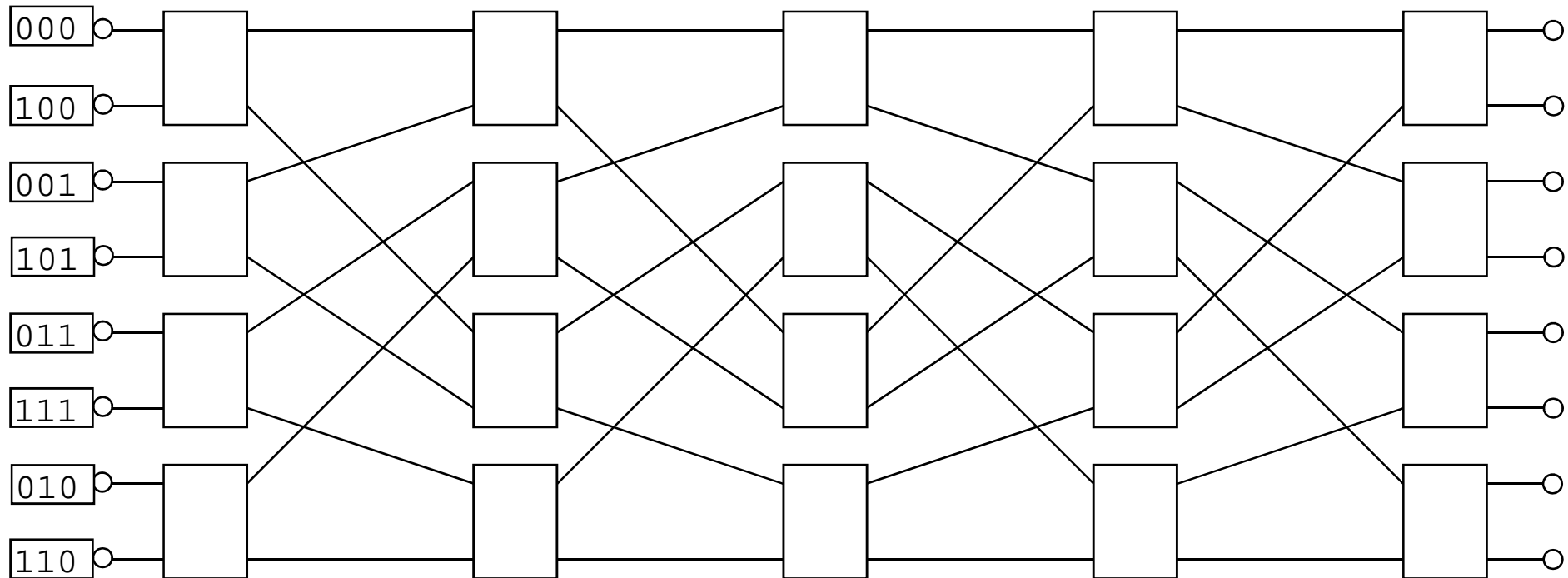
# Benešova síť (Beneš network)

## Samosměrovací algoritmus:

- Přepínač v  $i$ -tém, resp.  $j$ -tém stupni ( $1 \leq i \leq k$ , resp.  $k < j \leq 2k-1$ , kde  $k = \log_2 N$ ) síť čte  $i$ -tý bit, resp.  $(2k-j)$  tý bit směrovacího návěští (počínaje LSB směrem k MSB) a pokud se tento rovná:
  - 0 pak jdi na horní výstup
  - 1 pak jdi na dolní výstup

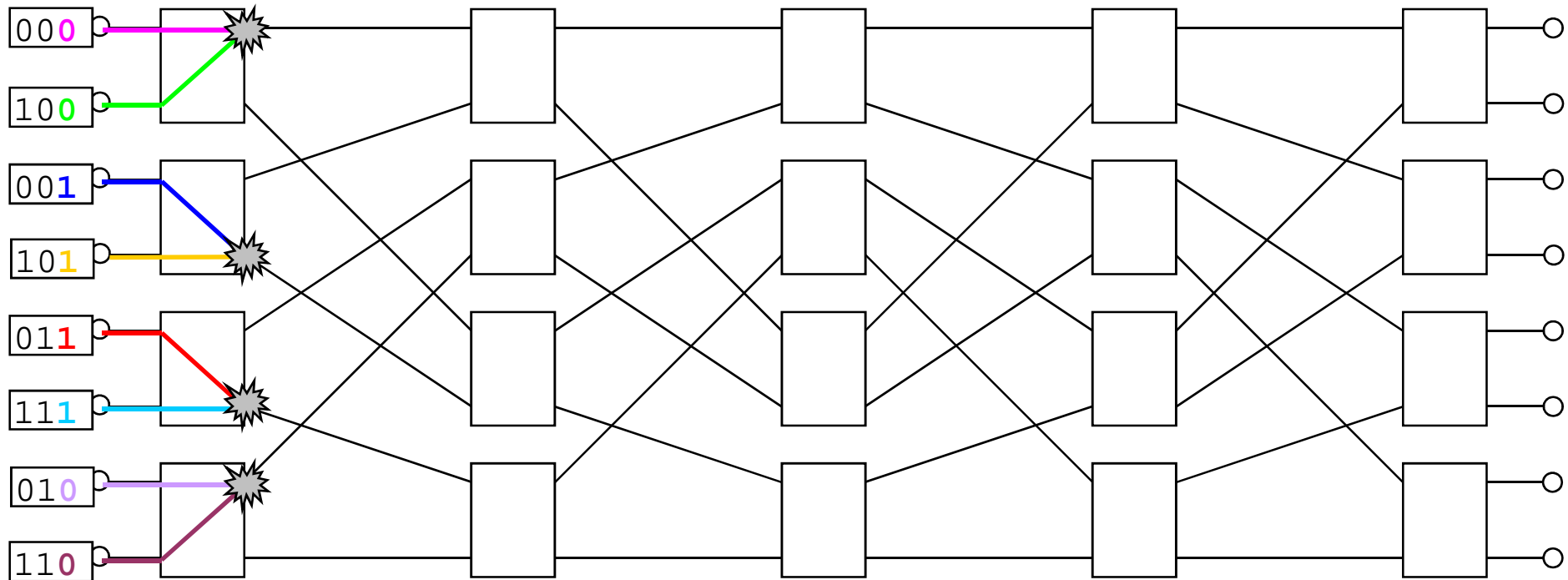


# Benešova síť (Beneš network)





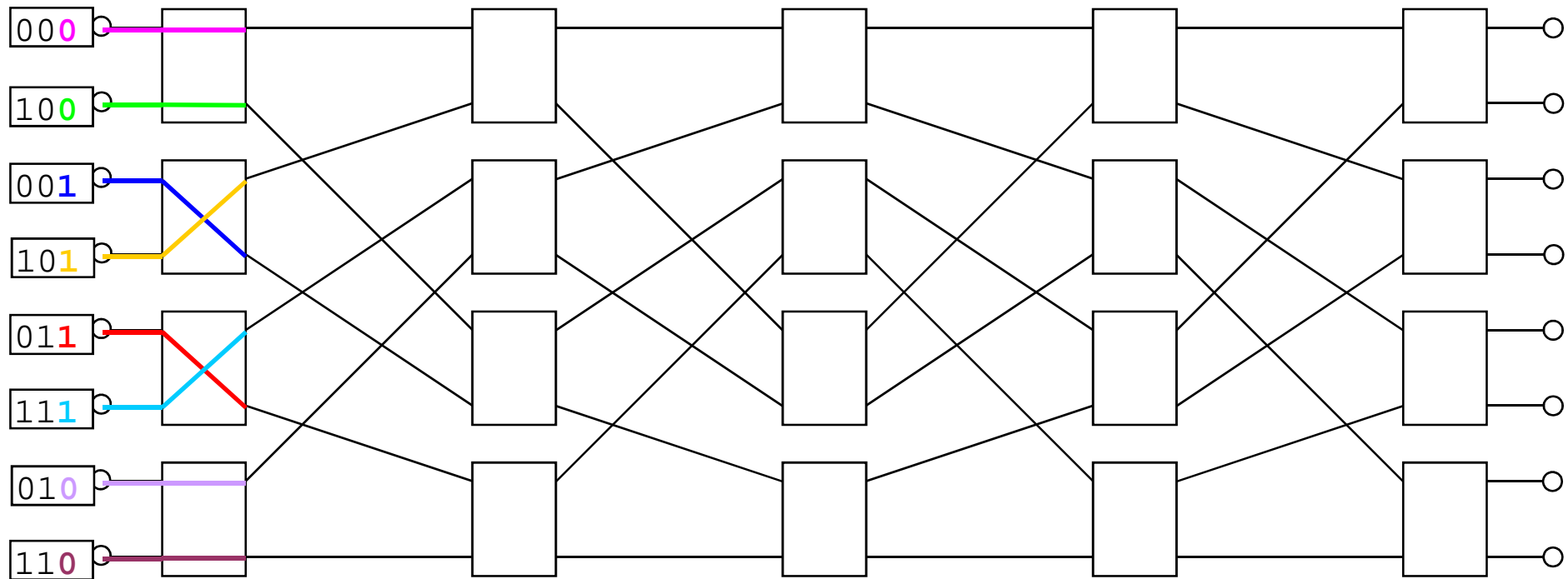
# Benešova síť (Beneš network)



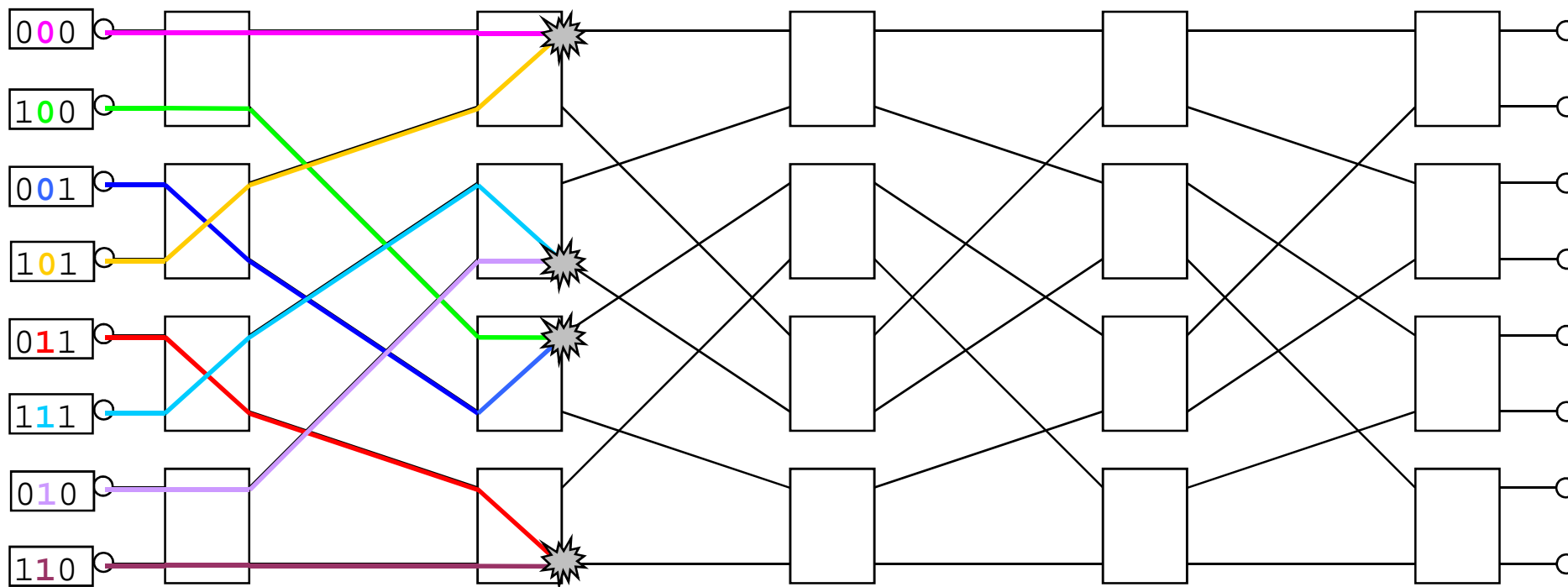
Co teď ???

Napr. Prednost má ten, kterého hodnota tagu (adresa určení) je menší...

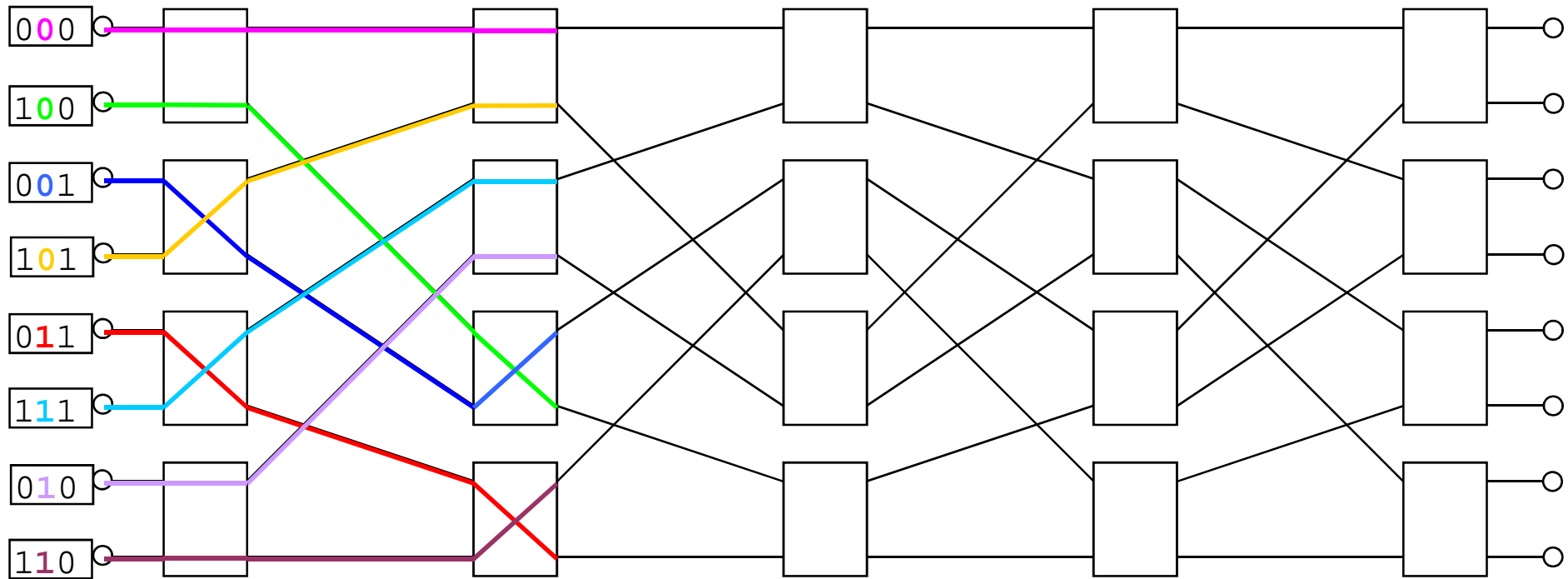
# Benešova síť (Beneš network)



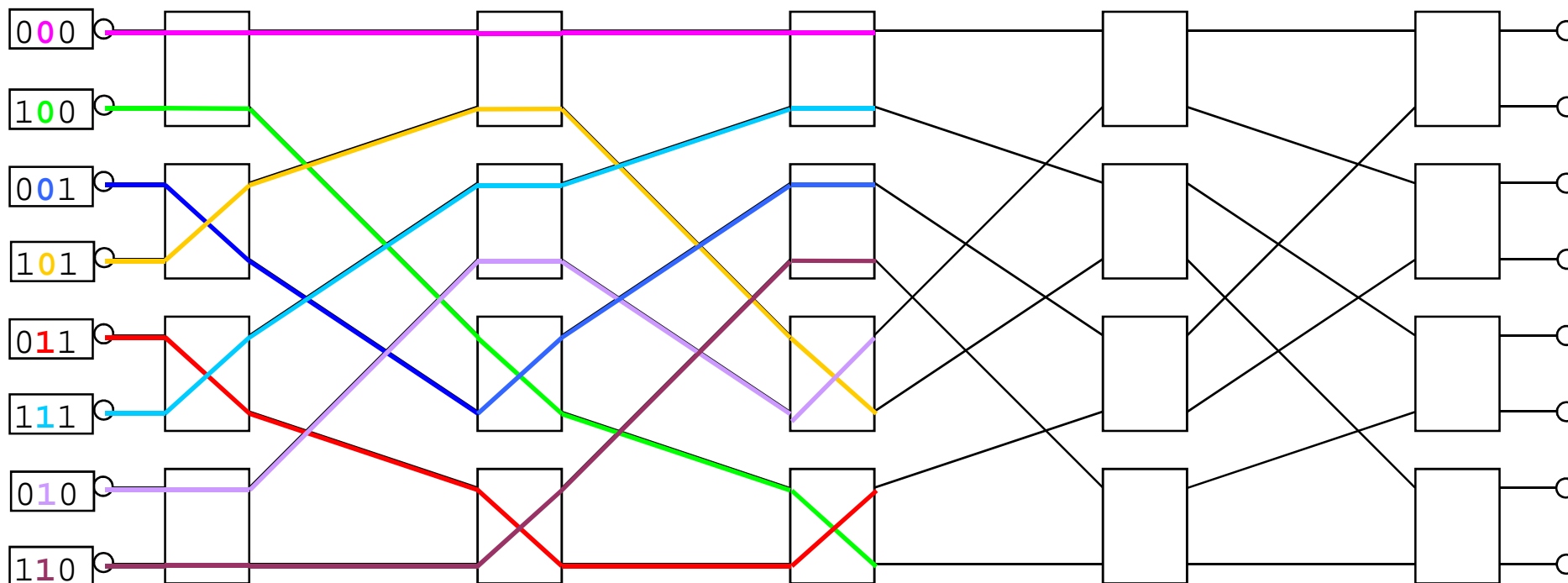
# Benešova síť (Beneš network)



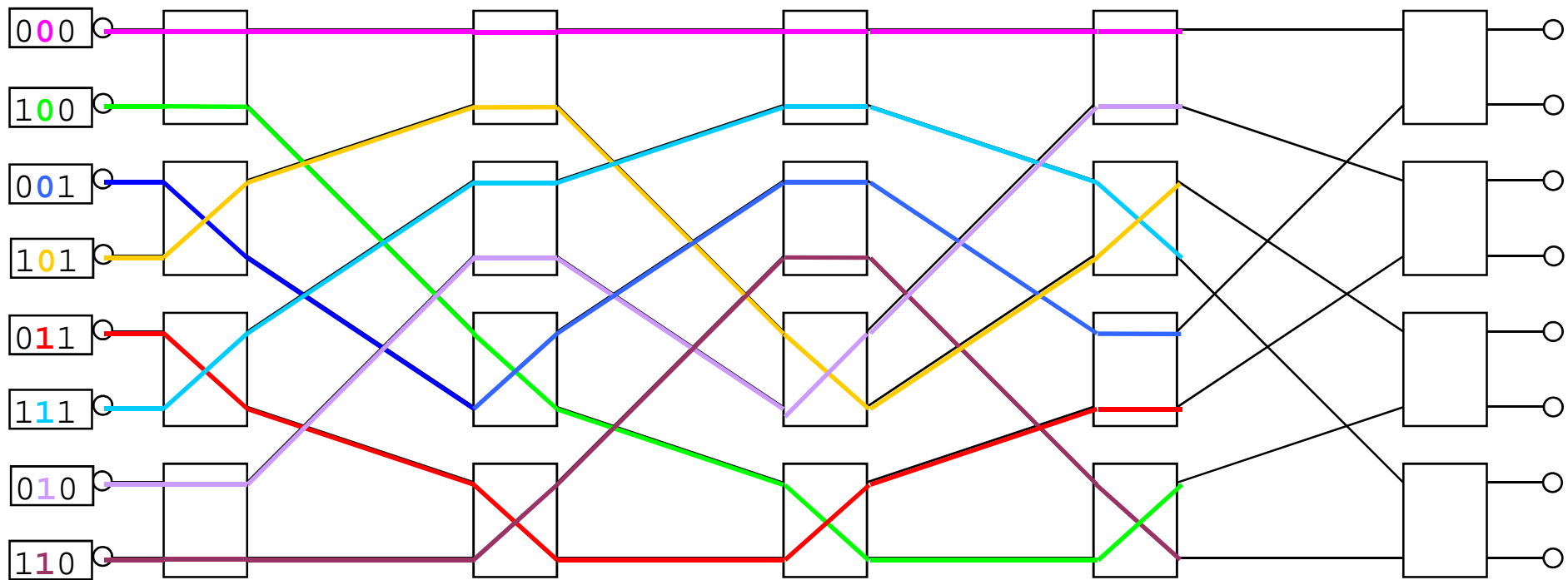
# Benešova síť (Beneš network)



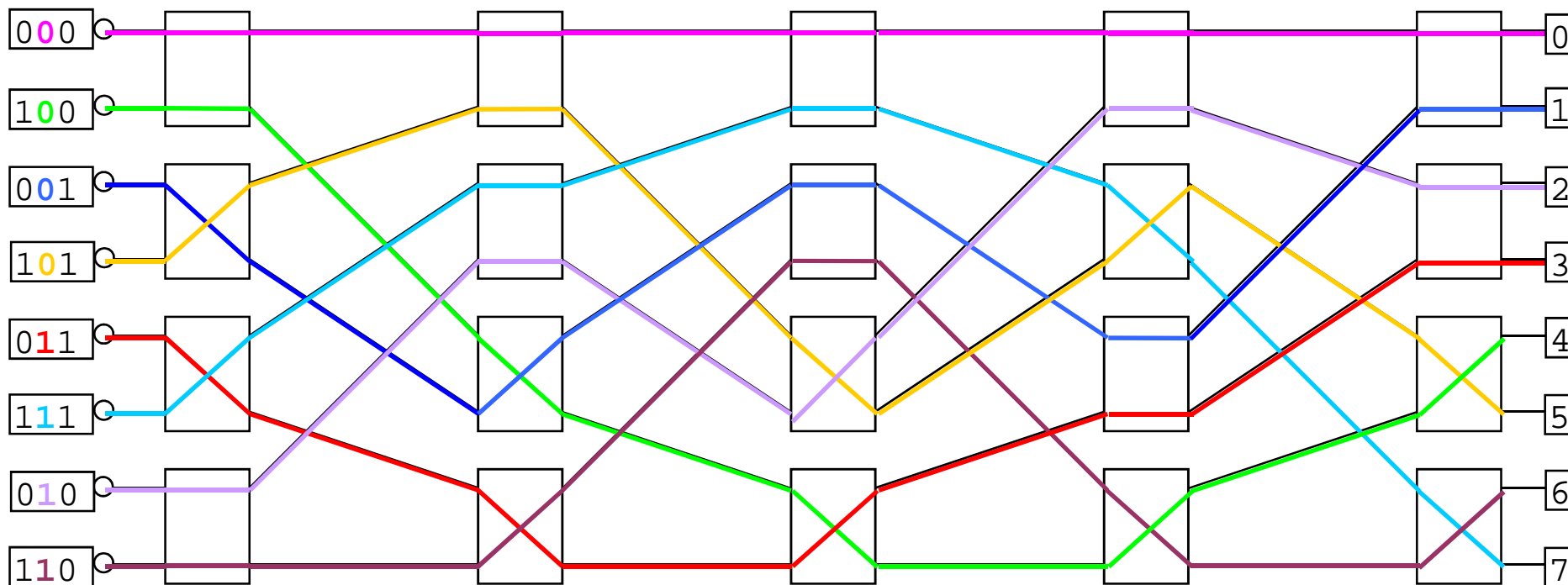
# Benešova síť (Beneš network)



# Benešova síť (Beneš network)

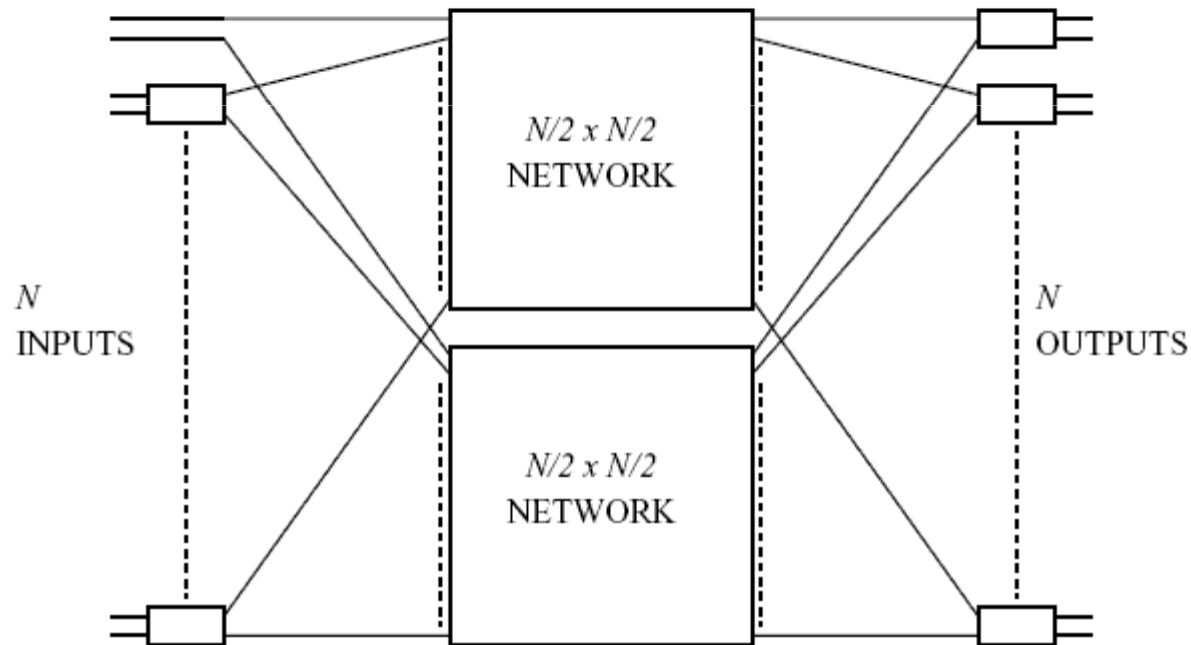


# Benešova síť (Beneš network)



## Waksmanova síť (Waksman network)

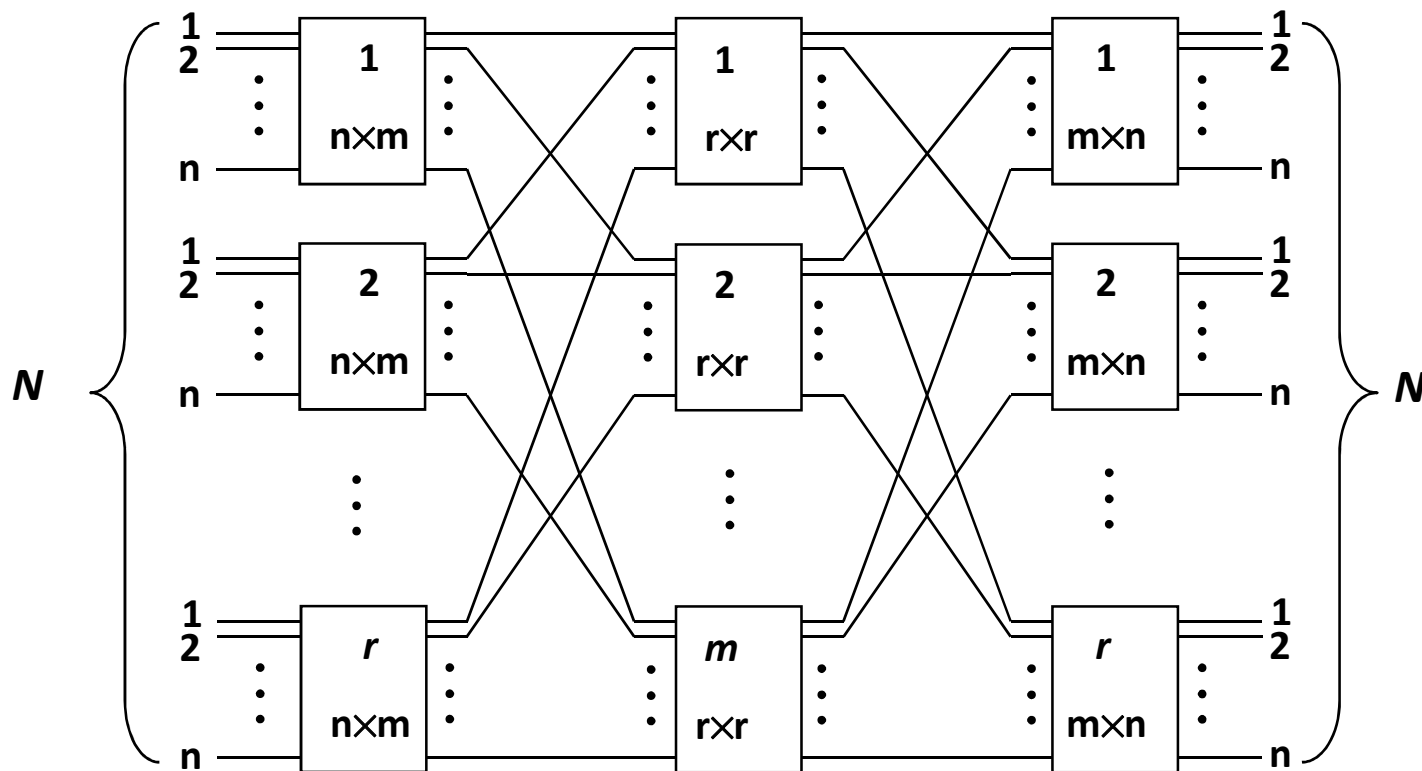
- Použitím modifikované verze Slepian-Duguidovho teorému může být z Benešovy sítě vynecháním jednoho vstupního nebo výstupního elementu vytvořena tzv. Waksmanova síť. Waksmanova síť používá méně přepínacích elementů než Benešova síť.





## Closova síť

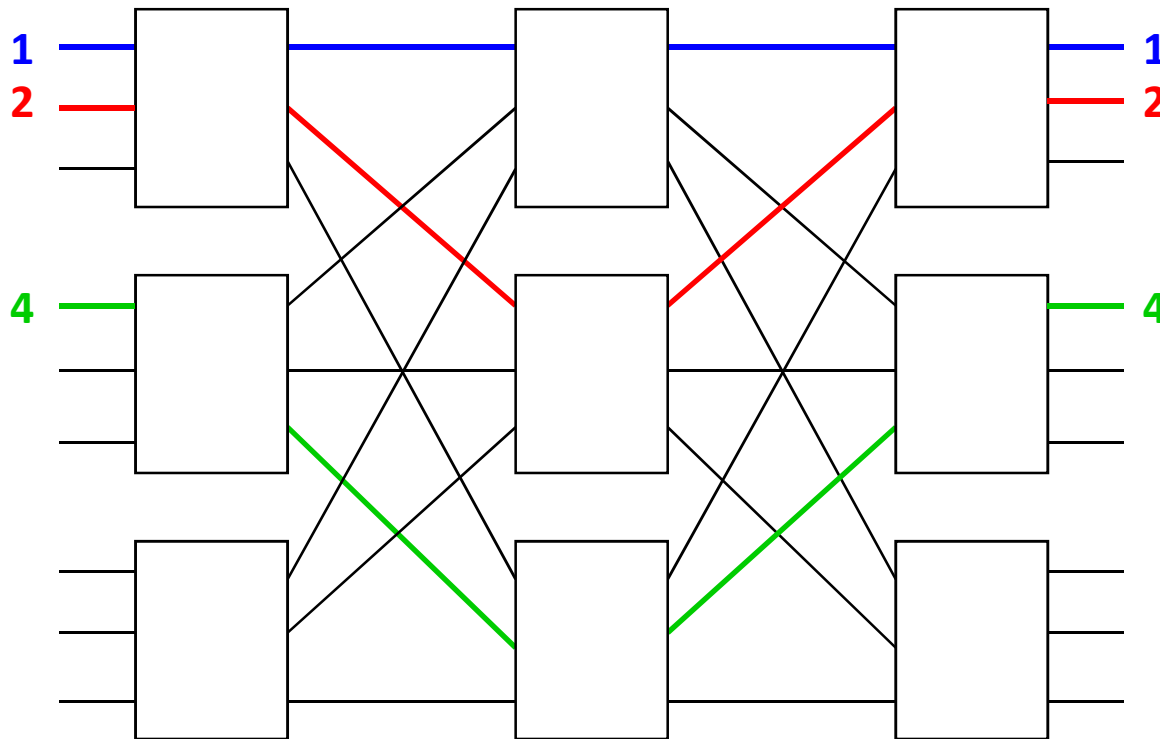
- Symetrická 3-stupňová Closova síť s  $N$  vstupními a  $N$  výstupními porty má  $r$  přepínacích modulů velikosti  $n \times m$  na prvním stupni,  $m$  přepínacích modulů velikosti  $r \times r$  na středním stupni a  $r$  přepínacích modulů velikosti  $m \times n$  na třetím (výstupním) stupni. Taková 3 stupňová síť je označována jako  $C(m, n, r)$ .



## Closova síť

Je Closova síť s  $m = n$  neblokující?

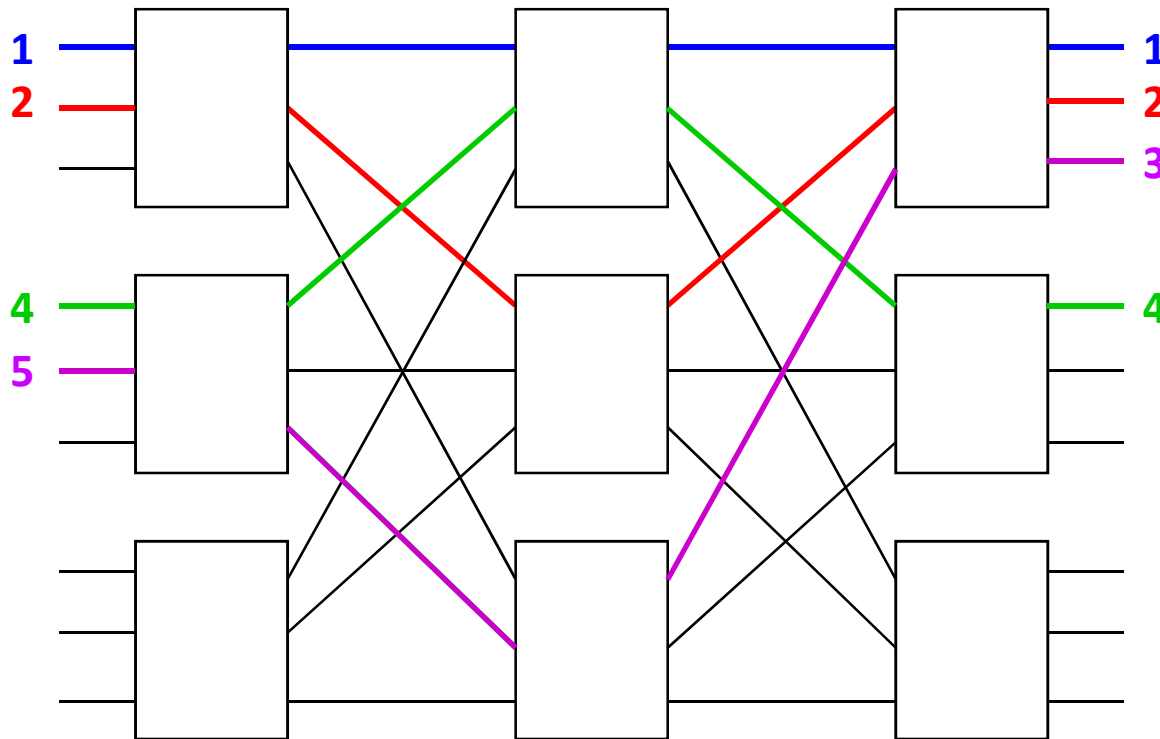
- **Příklad propojení v síti:** (1,1), (2,2), (4,4), (5,3), ...



## Closova síť

Je Closova síť s  $m = n$  neblokující?

- **Příklad propojení v síti:** (1,1), (2,2), (4,4), (5,3), ...

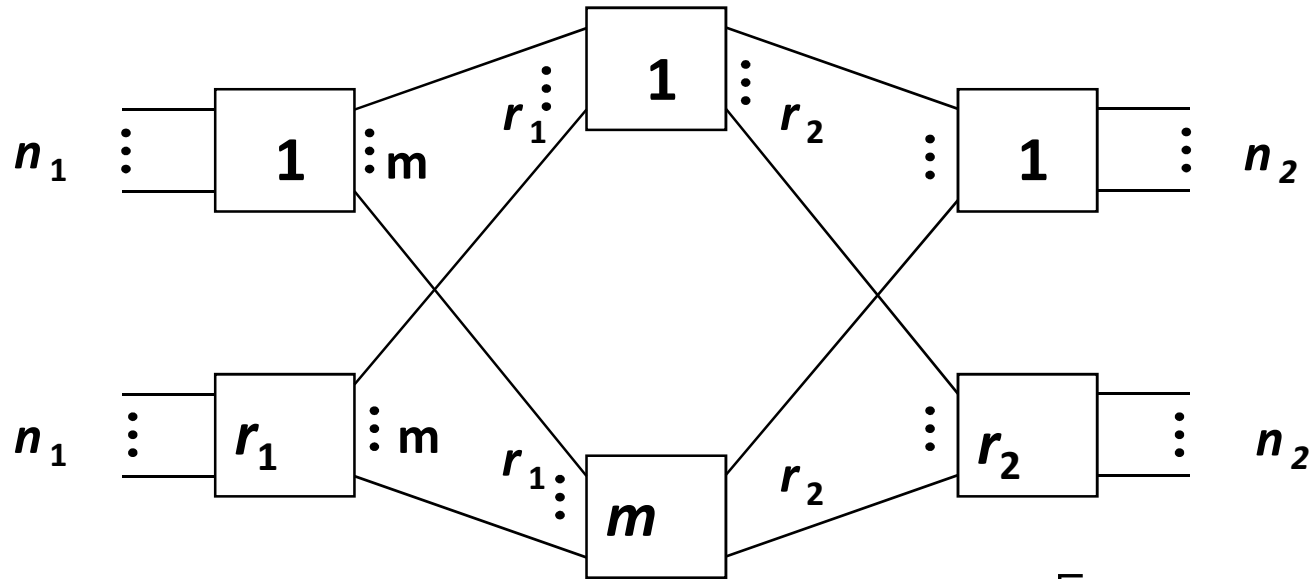


Řešením problému je přeuspořádat již existující spojení v síti

## Closova síť

- Propojovací schopnost Closovej sítě je závislá na parametrech  $n$ ,  $r$  a  $m$ . Pro dané  $n$ ,  $r$  a měnící se  $m$  je mnoho možností propojování. Neblokující operace 3-stupňové sítě můžeme dosáhnout více než jedním způsobem. Rozlišujeme 3 módy neblokujících Closových sítí:
  - **striktně (přísně) neblokující síť** (Strictly Nonblocking - SNB), podmínka:  $m \geq 2n - 1$ ,
  - **široce neblokující síť** (Wide-Sense Nonblocking - WSN),
  - **přeuspořadatelná neblokující síť** (Rearrangeably Nonblocking - RNB), podmínka:  $m \geq n$
- Obecně třístupňová Closova síť  $C(n_1, r_1, m, n_2, r_2)$  je 3-stupňová síť, jejíž první stupeň tvoří  $r_1$  přepínačů velikosti  $n_1 \times m$ , třetí stupeň má  $r_2$  přepínačů dimenze  $m \times n_2$ , a střední stupeň tvoří  $m$  přepínačů velikosti  $r_1 \times r_2$ . Pokud platí  $n_1 = n_2$ ,  $r_1 = r_2$  pak hovoříme o symetrické Closovej 3-stupňové síti.  $C(n, r, m, n, r)$  značíme  $C(m, n, r)$ .

## Closova síť



- Pro Closovu síť  $C(n_1, r_1, m, n_2, r_2)$  můžeme požadovaný stav sítě - zadaný permutací  $P$  přepsat do matice o  $r_1$  řádcích a  $r_2$  sloupcích

$$A = \begin{bmatrix} a_{1,1} & \cdots & a_{1,r_2} \\ \vdots & & \vdots \\ a_{r_1,1} & \cdots & a_{r_1,r_2} \end{bmatrix}$$

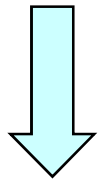
**Pro matici  $A$  platí:**  $\sum_{i=1}^{r_1} a_{ij} = n_2 \quad \text{pro } \forall j = \text{konšt.}$

$\sum_{j=1}^{r_2} a_{ij} = n_1 \quad \text{pro } \forall i = \text{konšt.}$

## Closova síť

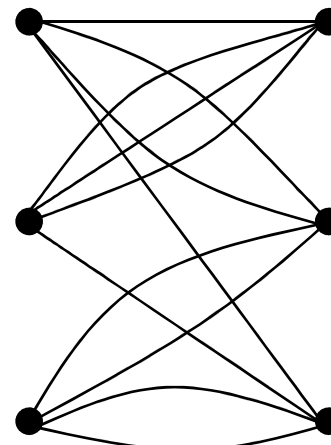
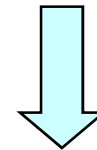
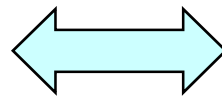
- Mějme permutaci  $P$  a síť typu  $C(4,4,3)$ :

$$P = \begin{pmatrix} 1 & 2 & 3 & 4 & \dots & 5 & 6 & 7 & 8 & \dots & 9 & 10 & 11 & 12 \\ 3 & 7 & 5 & 11 & \dots & 4 & 9 & 1 & 2 & \dots & 10 & 12 & 6 & 8 \end{pmatrix}$$



$$\begin{bmatrix} 1 & 2 & 1 \\ 3 & 0 & 1 \\ 0 & 2 & 2 \end{bmatrix}$$

propojovací matice



bipartitní graf

## Closova síť

Podle zvolené reprezentace vyplývají algoritmy:

- rozklad permutační matice
- rozklad bipartitního grafu

Rozklad propojovací matice:

- Rozkladem propojovací matice  $A$  na dílčí matice  $B_p$ ,  $p \in \{1, 2, \dots, q\}$ , také že platí  $A = B_1 + B_2 + \dots + B_p + \dots + B_q$  a zároveň pro  $\forall B_p$  platí
  - $b_{ij} \in \{0, 1\}$
  - $\sum_{i=1}^{r_1} b_{ij} = 1$  příp.  $\sum_{i=1}^{r_1} b_{ij} \leq 1$
  - $\sum_{j=1}^{r_2} b_{ij} = 1$  příp.  $\sum_{j=1}^{r_2} b_{ij} \leq 1$
- získáme konfiguraci přepínačů v druhém stupni sítě. Pro konkrétní rozklad je  $q!$  různých nastavení přepínačů v druhém stupni.

## Closova síť

- V našem případě:

$$\begin{bmatrix} 1 & 2 & 1 \\ 3 & 0 & 1 \\ 0 & 2 & 2 \end{bmatrix}$$



$$S_1 = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

$$S_2 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$S_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

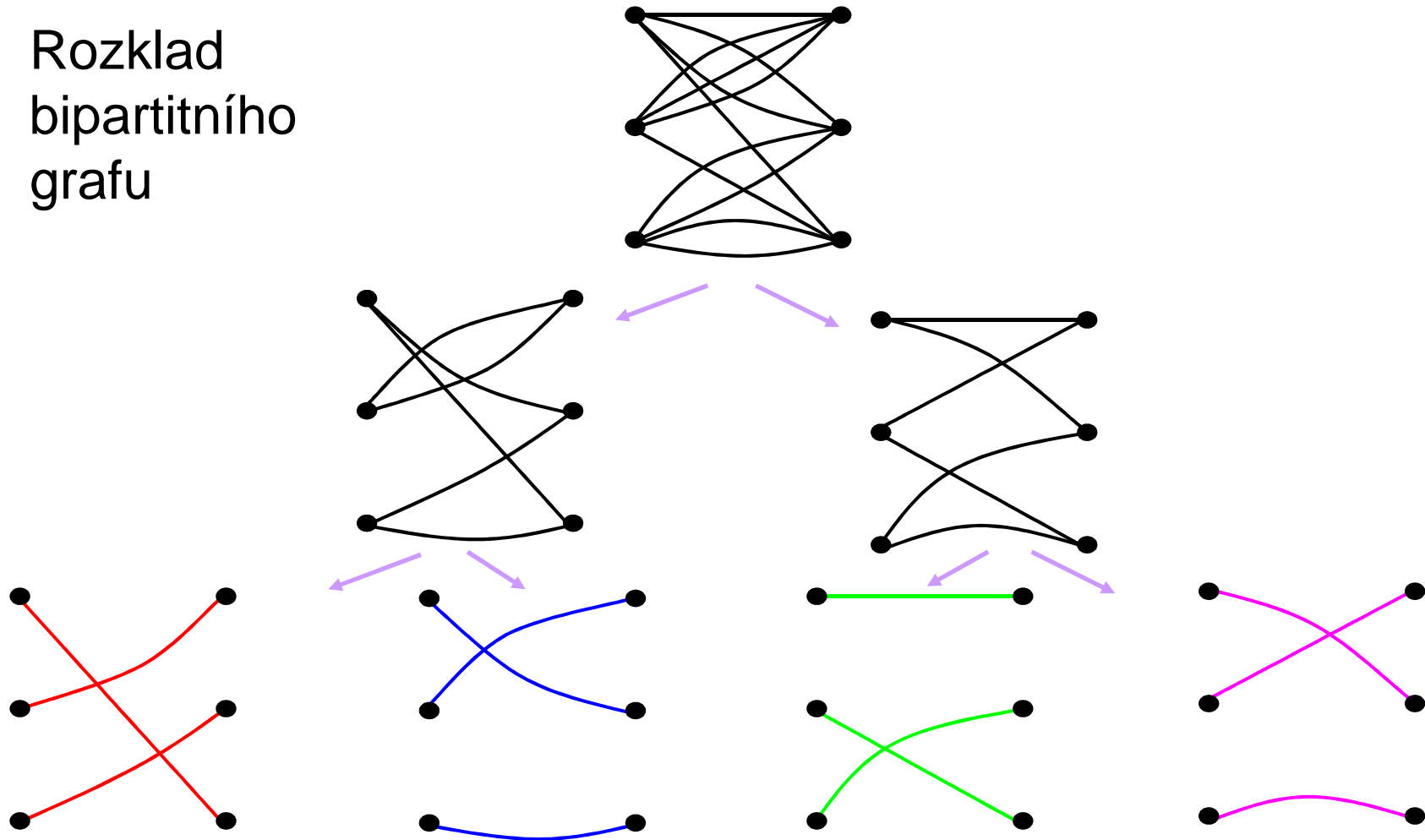
$$S_4 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\mathbf{A} = \mathbf{S1} + \mathbf{S2} + \mathbf{S3} + \mathbf{S4}$$



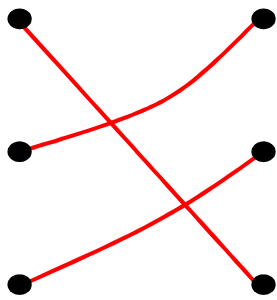
# Closova síť

- Rozklad bipartitního grafu

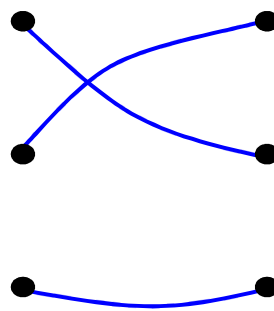


# Closova síť

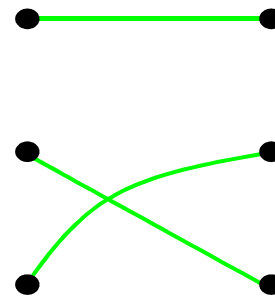
- Rozklad bipartitního grafu



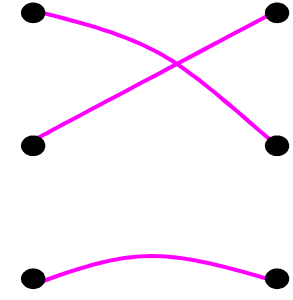
$$S_1 = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$



$$S_2 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$



$$S_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$



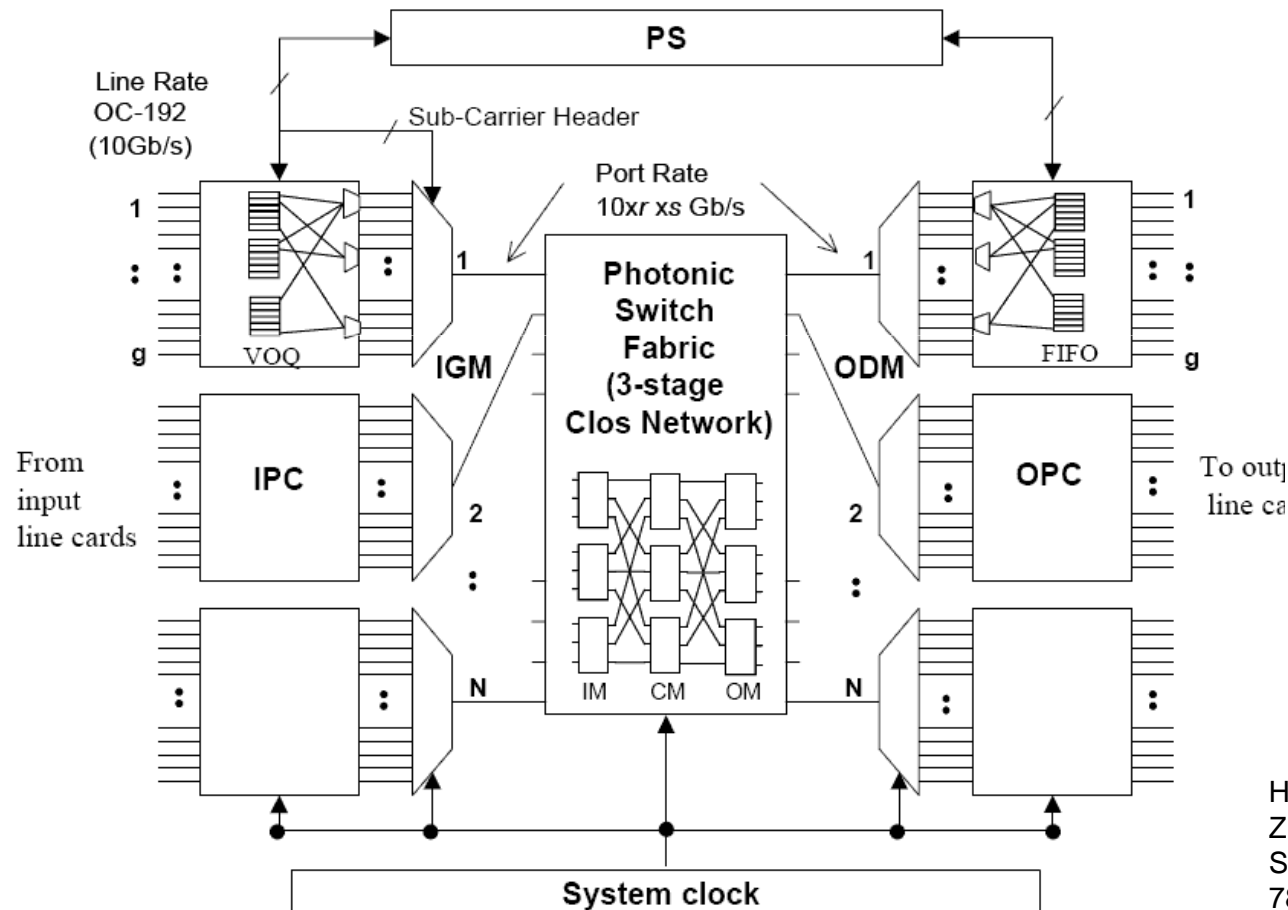
$$S_4 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

## Closova síť

- Matice vzniklé rozkladem (ať už grafu nebo propojovací matice) umožňují **přímo** nastavit přepínače ve středním stupni sítě.
- Pokud známe nastavení přepínačů středního stupně není náročné nastavit přepínače v celé síti.
- Pozn: Benešova síť je speciální případ Closovej sítě. Proto pokud chceme zajistit bezkonfliktnost Benešovy sítě můžeme využít výše uvedené algoritmy pro Closovu síť.

# Closova síť

- Například architektura peta bytového paketového optického přepínače P3S dimenze 6400x6400 je založena na Closovej 3-stupňové síti dosahuje celkově 1.024 petabit / s (160 Gbit / s na každý port)

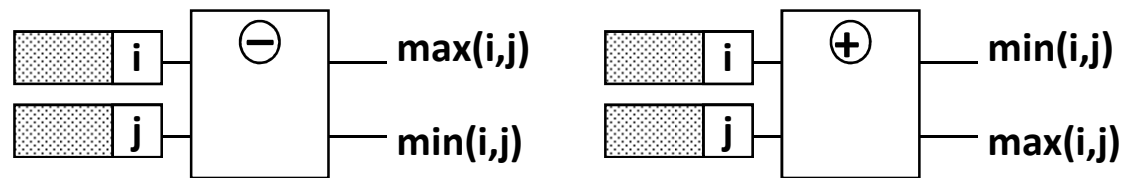


IPC: Input Port Interface Card  
 IGM: Input Grooming Module  
 ODM: Output Demultiplexing Module  
 OPC: Output Port Interface Card  
 PS: Packet Scheduler  
 VOQ: Virtual Output Queue  
 r: Cell Number / s: Speedup  
 g: Input Line Number  
 PSF: Photonic Switching Fabric  
 IM: Input Module  
 CM: Central Module  
 OM: Output Module

H. Jonathan Chao, Kung-Li Deng, and Zhigang Jing: A Petabit Photonic Packet Switch (P<sup>3</sup>S), IEEE INFOCOM 2003, 0-7803-7753-2/03

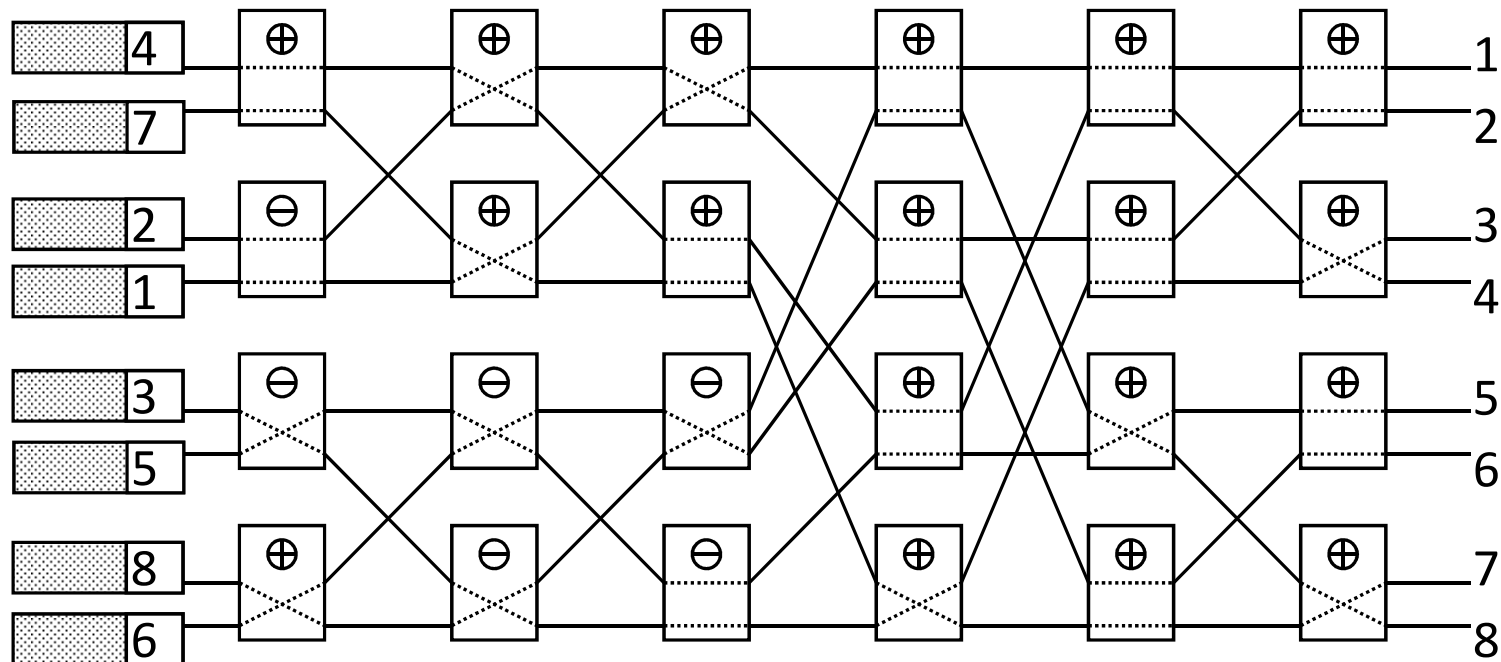
## Batcherova síť

- Batcherova, nebo třídící (sorting) síť řadí vstupující pakety podle jejich adresy výstupu a to od nejmenších adres k největším.
- Batcherova třídící síť je složena z elementů:



- Vhodným seřazením takových elementů dostaneme Batcherovu síť. Celá třídící síť je složena z posloupnosti **bitonických** řadičů (bitonic sorters), které řadí své výstupy sestupně, nebo vzestupně. V našem případě jsou to třídiče velikosti  $2 \times 2$ . Pokud je na vstupu elementu jeden jen paket a tedy jen jedna adresa, uvažuje se jak by měla nižší hodnotu.

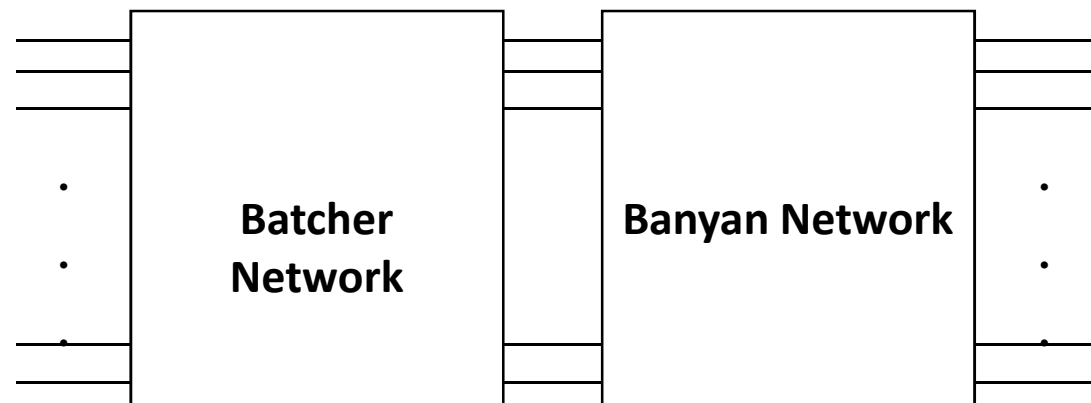
## Batcherova síť



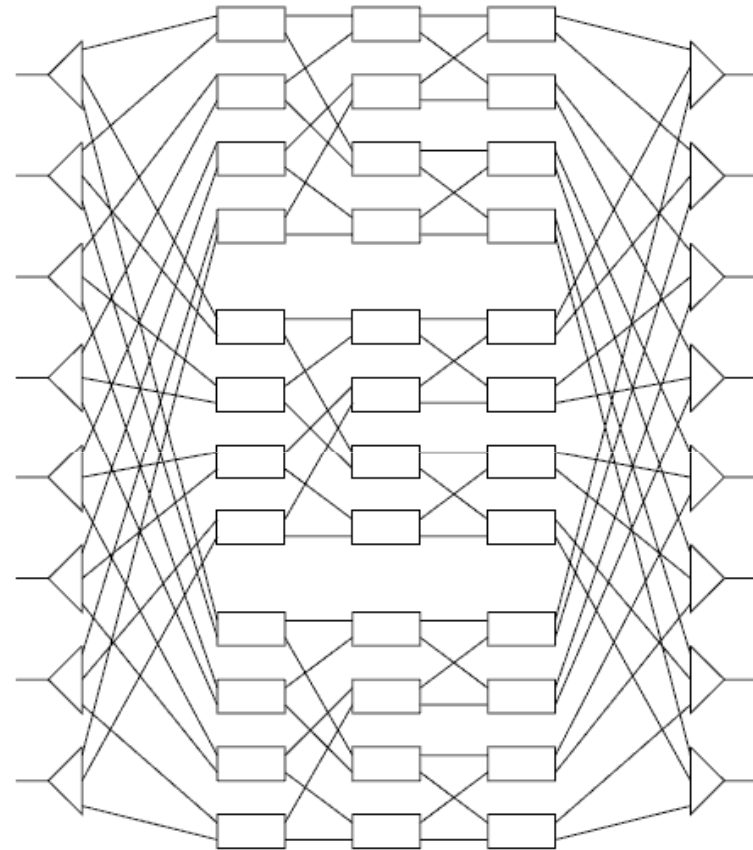
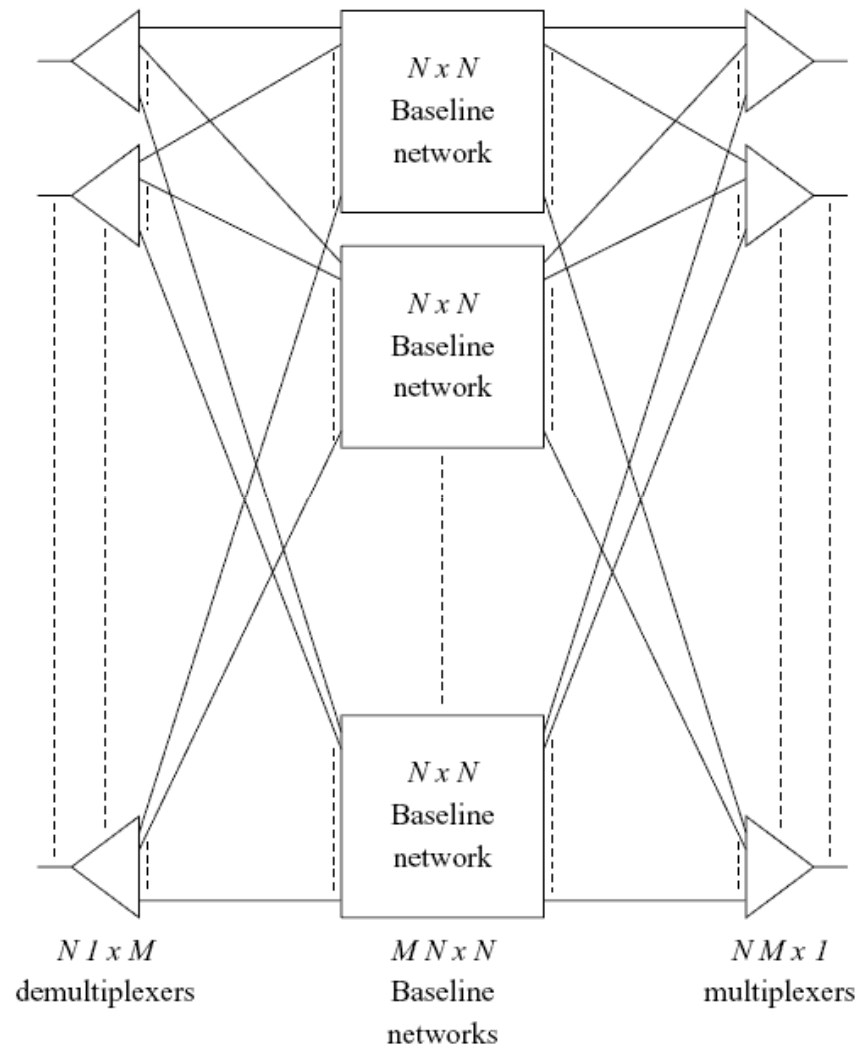
- Na výstupu Batcherovej síťe jsou pakety seřazeny podle vzestupných adres, ale nedosáhnou správný výstup podle své adresy určení (pokud nejsou obsazeny všechny vstupy). Proto se Batcherova síť dále kombinuje s banyan sítí.

## Batcher-Banyan síť

- Banyan síť zapojena za Batcherovou sítí má samosmerovacie vlastnosti a dopraví tak pakety ke správnému výstupu ze sítě. Zařazení třídící sítě před banyan sítí vyloučí blokování HOL pokud předpokládáme, že žádné dva pakety nejsou směřovány k témuž výstupu. Vzhledem k tomu, že pakety jsou roztríděna v třídící síti, nevznikne ani vnitřní blokování v banyan síti. V případě pokud existuje možnost stejných výstupních adres na vstupu sítě, je třeba použít vyrovnávací paměti.



# Paralelní Baseline síť





## Paralelní Baseline síť

- Na vstupu v výstupu sítě paralelní baseline sítě (někdy také označované jako **Multi-Log<sub>2</sub>N** network) je realizována funkce expanze a koncentrace. Na vstupu sítě jsou spojovací elementy  $1 \times m$ , které expandují provoz a na výstupu elementy  $m \times 1$ , které provoz koncentrují.
- Pokud nepočítáme vstupní a výstupní stupně (expanze a koncentrace), počet stupňů v síti je  $n = \log_2 N$ . Z teorie grafů je možné dokázat, že pokud je počet paralelních baseline podsítí

$$m \geq 2^{(n/2)}$$

pak paralelní baseline síť je rekonfigurovatelná bez blokády. To zároveň znamená, že pro optimální paralelní baseline síť potřebujeme více spojovacích elementů než pro sériovou baseline síť (Benešovu síť). Na druhou stranu, čas přechodu paralelními baseline sítí je menší než při sériové baseline.

## Cantor network

- Podobně jako Closova síť, i Cantorova síť patří mezi striktně neblokující. Cantorova síť  $N \times N$  může být vytvořena z  $\log_2 N$  Benešových sítí,  $N$  demultiplexů a  $N$  multiplexorů.

