

ESE680-002 (ESE534): Computer Organization

Day 16: March 14, 2007
Interconnect 4: Switching



Penn ESE680-002 Spring2007 -- DeHon

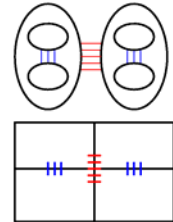
Previously

- Used Rent's Rule characterization to understand wire growth

$$IO = c N^p$$

- Top bisections will be $\Omega(N^p)$
- 2D wiring area

$$\Omega(N^p) \times \Omega(N^p) = \Omega(N^{2p})$$



Penn ESE680-002 Spring2007 -- DeHon

We Know

- How we avoid $O(N^2)$ wire growth for "typical" designs
- How to characterize locality
- How we might exploit that locality to reduce wire growth
- Wire growth implied by a characterized design

3

Penn ESE680-002 Spring2007 -- DeHon

Today

- Switching
 - Implications
 - Options

4

Penn ESE680-002 Spring2007 -- DeHon

Switching:

How can we use the locality captured by Rent's Rule to reduce switching requirements? (How much?)

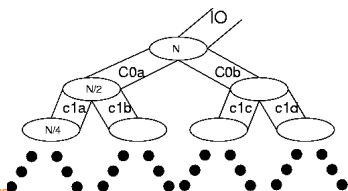
5

Penn ESE680-002 Spring2007 -- DeHon

Observation

- Locality that saved us wiring, also saves us switching

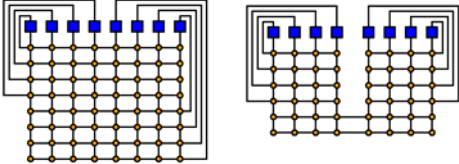
$$IO = c N^p$$



Penn ESE680-002 Spring2007 -- DeHon

Consider

- Crossbar case to exploit wiring:
 - split into two halves, connect with limited wires
 - $N/2 \times N/2$ crossbar each half
 - $N/2 \times (N/2)^p$ connect to bisection wires
 - $2(N^2/4) + 2(N/2)^{p+1} < N^2$

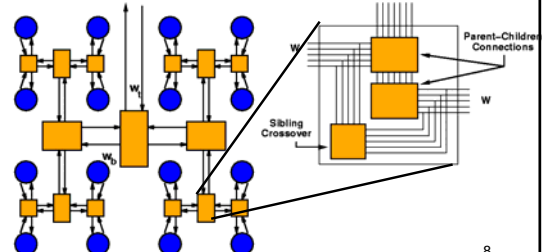


Penn ESE680-002 Spring2007 -- DeHon

7

Recurse

- Repeat at each level
 - form tree

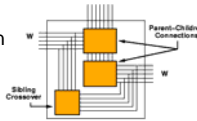


Penn ESE680-002 Spring2007 -- DeHon

8

Result

- If use crossbar at each tree node
 - $O(N^{2p})$ wiring area
 - for $p > 0.5$, direct from bisection
 - $O(N^{2p})$ switches
 - top switch box is $O(N^{2p})$
 - $2 W_{top} \times W_{bot} + (W_{bot})^2$
 - $2 (N^p \times (N/2)^p) + (N/2)^{2p}$
 - $N^{2p}(1/2^p + 1/2^{2p})$

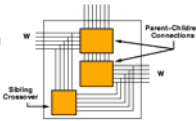


Penn ESE680-002 Spring2007 -- DeHon

9

Result

- If use crossbar at each tree node
 - $O(N^{2p})$ wiring area
 - for $p > 0.5$, direct from bisection
 - $O(N^{2p})$ switches
 - top switch box is $O(N^{2p})$
 - $N^{2p}(1/2^p + 1/2^{2p})$
 - switches at one level down is
 - $2 \times (N/2)^{2p}(1/2^p + 1/2^{2p})$
 - $2 \times (1/2^p)^2 \times (N^{2p}(1/2^p + 1/2^{2p}))$
 - $2 \times (1/2^p)^2 \times \text{previous level}$

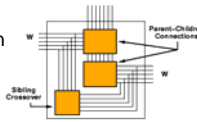


Penn ESE680-002 Spring2007 -- DeHon

10

Result

- If use crossbar at each tree node
 - $O(N^{2p})$ wiring area
 - for $p > 0.5$, direct from bisection
 - $O(N^{2p})$ switches
 - top switch box is $O(N^{2p})$
 - $N^{2p}(1/2^p + 1/2^{2p})$
 - switches at one level down is
 - $2 \times (1/2^p)^2 \times \text{previous level}$
 - $(2/2^{2p}) = 2^{(1-2p)}$
 - coefficient < 1 for $p > 0.5$

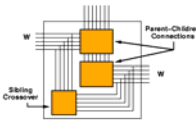


Penn ESE680-002 Spring2007 -- DeHon

11

Result

- If use crossbar at each tree node
 - $O(N^{2p})$ switches
 - top switch box is $O(N^{2p})$
 - switches at one level down is
 - $2^{(1-2p)} \times \text{previous level}$
 - Total switches:
 - $N^{2p} \times (1 + 2^{(1-2p)} + 2^{2(1-2p)} + 2^{3(1-2p)} + \dots)$
 - get geometric series; sums to $O(1)$
 - $N^{2p} \times (1 / (1 - 2^{(1-2p)}))$
 - $= 2^{(2p-1)} / (2^{(2p-1)} - 1) \times N^{2p}$



Penn ESE680-002 Spring2007 -- DeHon

12

Good News

- Good news
 - asymptotically optimal
 - Even without switches, area $O(N^{2p})$
 - so adding $O(N^{2p})$ switches not change

Penn ESE680-002 Spring2007 -- DeHon

13

Bad News

- Switches area \gg wire crossing area
 - Consider 8λ wire pitch \Rightarrow crossing $64\lambda^2$
 - Typical (passive) switch $\Rightarrow 2500\lambda^2$
 - Passive only: 40x area difference
 - worse once rebuffer or latch signals.
- ...and switches limited to substrate
 - whereas can use additional metal layers for wiring area

Penn ESE680-002 Spring2007 -- DeHon

14

Additional Structure

- This motivates us to look beyond crossbars
 - can depopulate crossbars on up-down connection without loss of functionality?

Penn ESE680-002 Spring2007 -- DeHon

15

Can we do better?

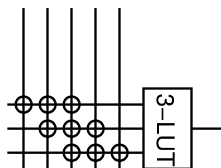
- Crossbar too powerful?
 - Does the specific down channel matter?
- What do we want to do?
 - Connect to *any* channel on lower level
 - Choose a subset of wires from upper level
 - order not important

Penn ESE680-002 Spring2007 -- DeHon

16

N choose K

- Exploit freedom to depopulate switchbox
- Can do with:
 - $K \times (N-K+1)$ switches
 - Vs. $K \times N$
 - Save $\sim K^2$

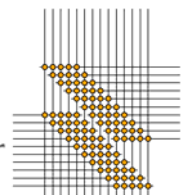


Penn ESE680-002 Spring2007 -- DeHon

17

N-choose-M

- Up-down connections
 - only require concentration
 - choose M things out of N
 - i.e. **order** of subset irrelevant
- Consequent:
 - can save a constant factor $\sim 2^p/(2^p-1)$
 - $(N/2)^p \times N^p$ vs $(N^p - (N/2)^{p+1})/(N/2)^p$
 - $p=2/3 \rightarrow 2^p/(2^p-1) \approx 2.7$
- Similary, Left-Right
 - order not important \Rightarrow reduces switches



Penn ESE680-002 Spring2007 -- DeHon

18

Multistage Switching

Penn ESE680-002 Spring2007 -- DeHon

19

Multistage Switching

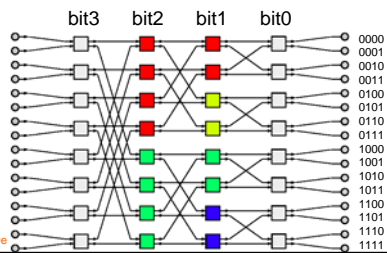
- We can route any **permutation** w/ less switches than a crossbar
- If we allow switching in stages
 - Trade increase in switches in path
 - For decrease in total switches

Penn ESE680-002 Spring2007 -- DeHon

20

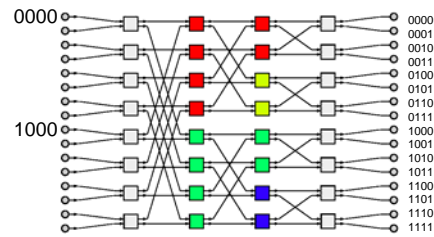
Butterfly

- Log stages
- Resolve one bit per stage



Penn ESE680-002 Spring2007 -- De

What can a Butterfly Route?

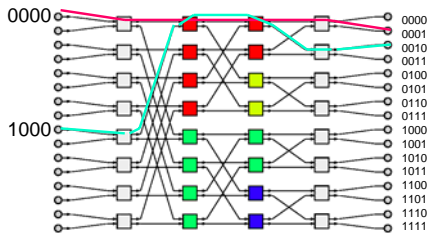


- 0000 → 0001
- 1000 → 0010

Penn ESE680-002 Spring2007 -- DeHon

22

What can a Butterfly Route?



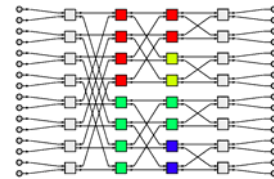
- 0000 → 0001
- 1000 → 0010

Penn ESE680-002 Spring2007 -- DeHon

23

Butterfly Routing

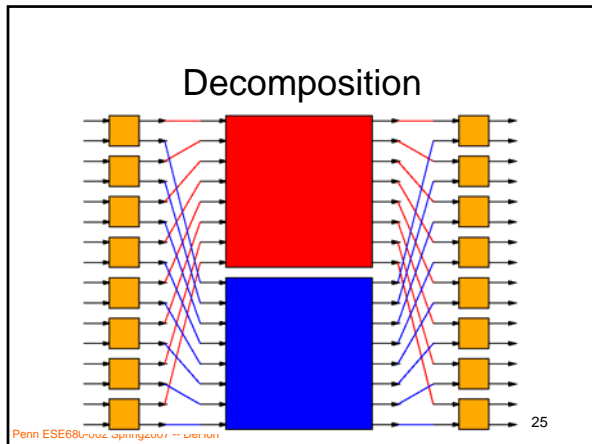
- **Cannot** route all permutations
 - Get internal blocking



- What required for non-blocking network?

Penn ESE680-002 Spring2007 -- DeHon

24



Decomposed Routing

- Pick a link to route.
- Route to destination over **red** network
- At destination,
 - What can we say about the link which shares the final stage switch with this one?
 - What can we do with link?
- Route that link
 - What constraint does this impose?
 - So what do we do?

26

Penn ESE680-002 Spring2007 -- DeHon

Decomposition

- Switches: $N/2 \times 2 \times 4 + (N/2)^2 < N^2$

27

Penn ESE680-002 Spring2007 -- DeHon

Recurse

If it works once, try it again...

28

Penn ESE680-002 Spring2007 -- DeHon

Result: Beneš Network

- $2\log_2(N)-1$ stages (switches in path)
- Made of $N/2$ 2×2 switchpoints
 - (4 switches)
- $4N \times \log_2(N)$ total switches
- Compute route in $O(N \log(N))$ time

29

Penn ESE680-002 Spring2007 -- DeHon

Beneš Network Wiring

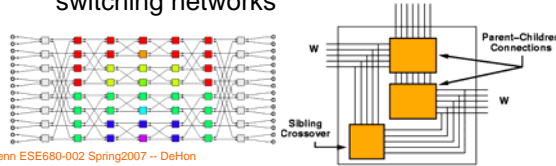
- Bisection: N
- Wiring $\rightarrow O(N^2)$ area (fixed wire layers)

30

Penn ESE680-002 Spring2007 -- DeHon

Beneš Switching

- Beneš reduced switches
 - N^2 to $N(\log(N))$
 - using multistage network
- Replace crossbars in tree with Beneš switching networks



Penn ESE680-002 Spring2007 -- DeHon

Beneš Switching

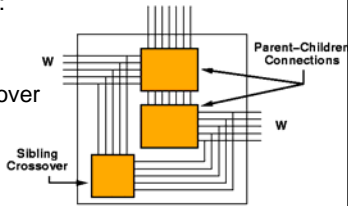
- Implication of Beneš Switching
 - still require $O(W^2)$ wiring per tree node
 - or a total of $O(N^{2p})$ wiring
 - now $O(W \log(W))$ switches per tree node
 - converges to $O(N)$ total switches!
 - $O(\log^2(N))$ switches in path across network
 - strictly speaking, dominated by wire delay $\sim O(N^p)$
 - but constants make of little practical interest except for very large networks ☹

Penn ESE680-002 Spring2007 -- DeHon

32

Better yet...

- Believe do not need Beneš on the up paths
- Single switch on up path
- Beneš for crossover
- Switches in path:
 - $\log(N)$ up
 - $+$ $\log(N)$ down
 - $+$ $2\log(N)$ crossover
 - $= 4 \log(N)$
 - $= O(\log(N))$



Penn ESE680-002 Spring2007 -- DeHon

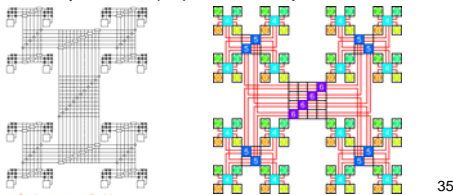
Linear Switch Population

Penn ESE680-002 Spring2007 -- DeHon

34

Linear Switch Population

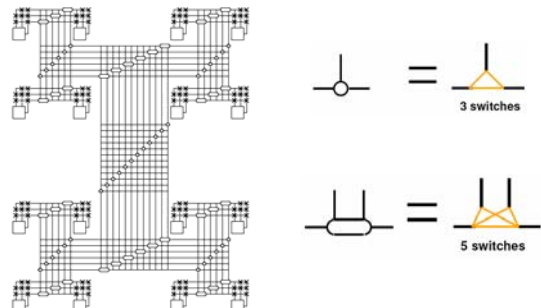
- Can further reduce switches
 - connect each lower channel to $O(1)$ channels in each tree node
 - end up with $O(W)$ switches per tree node



Penn ESE680-002 Spring2007 -- DeHon

35

Linear Switch ($p=0.5$)

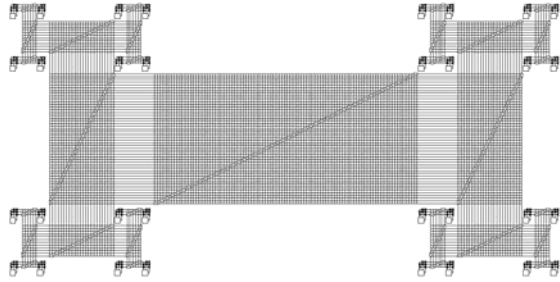


Penn ESE680-002 Spring2007 -- DeHon

36

Linear Population and Beneš

- Top-level crossover of $p=1$ is Beneš switching

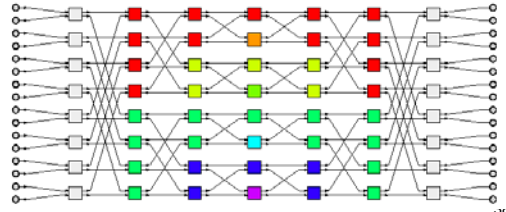


Penn ESE680-002 Spring2007 -- DeHon

37

Beneš Compare

- Can permute stage switches so local shuffles on outside and big shuffle in middle



Penn ESE680-002 Spring2007 -- DeHon

38

Linear Consequences: Good News

- Linear Switches
 - $O(\log(N))$ switches in path
 - $O(N^{2p})$ wire area
 - $O(N)$ switches
- More practical than Beneš crossover case

Penn ESE680-002 Spring2007 -- DeHon

39

Linear Consequences: Bad News

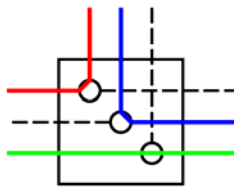
- Lacks guarantee can use all wires
 - as shown, at least mapping ratio > 1
 - likely cases where even **constant** not suffice
 - expect no worse than logarithmic
- Finding Routes is harder
 - no longer linear time, deterministic
 - **open** as to exactly how hard

Penn ESE680-002 Spring2007 -- DeHon

40

Mapping Ratio

- Mapping ratio says
 - if I have W channels
 - may only be able to use W/MR wires
 - for a particular design's connection pattern
 - to accommodate any design
 - for all channels



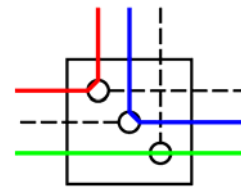
Penn ESE680-002 Spring2007 -- DeHon

$$\text{physical wires} \geq MR \times \text{logical}$$

41

Mapping Ratio

- Example:
 - Shows $MR=3/2$
 - For Linear Population, 1:1 switchbox



Penn ESE680-002 Spring2007 -- DeHon

42

Area Comparison

Both:
 $p=0.67$
 $N=1024$

M-choose-N
perfect map
Linear
MR=2

Penn ESE680-002 Spring2007 -- DeHon

43

Area Comparison

M-choose-N
perfect map
Linear
MR=2

- Since
 - switch \gg wire
- may be able to tolerate $MR > 1$
- reduces switches
 - net area savings
- Empirical:
 - Never seen greater than 1.5

Penn ESE680-002 Spring2007 -- DeHon

44

Expanders

Penn ESE680-002 Spring2007 -- DeHon

45

Expander Theory

□ (α, β) -expansion

- Any group of size $k = \alpha N$ will expand connect to a group of size $\beta k = \beta \alpha N$ in each logical direction

[Arora, Leighton, Maggs
SIAM Journal of Comp. v25n3p600 1996]

Penn ESE680-002 Spring2007 -- DeHon

46

Expander Idea

- **IF** we can achieve expansion
 - Can guarantee non-blocking at each stage
- *i.e.*
 - Guarantee use less than αN
 - Guarantee connections to more stuff in next level
 - Since $\beta \alpha N > \alpha N$ available in next level
 - Guaranteed to be an available switch

Penn ESE680-002 Spring2007 -- DeHon

47

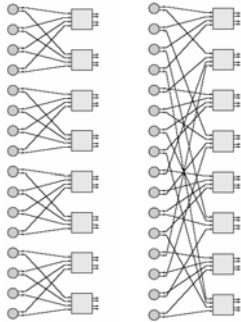
Dilated Switches

- Have multiple outputs per logical direction
 - **Dilation:** number of outputs per direction
 - *E.g.* radix 2 switch w/ 4 outputs
 - 2 per direction
 - Dilation 2

Penn ESE680-002 Spring2007 -- DeHon

48

Dilated Switches allow Expansion



- On Right
 - Any pair of nodes connects to 3 switches
- Strictly speaking must have $d > 2$ for expansion

Penn ESE680-002 Spring2007 -- DeHon

49

Random Wiring

- Random, dilated wiring for butterfly can achieve

$$d > \beta + 1 + \frac{\beta + 1 + \ln 2\beta}{\ln(\frac{1}{2\alpha\beta})}$$

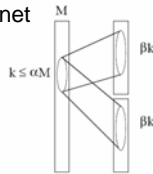
$$2d > 2\beta + 1 + \frac{2\beta + 1 + \ln 2\beta}{\ln(\frac{1}{2\alpha\beta})}$$

- For tree... $2 \rightarrow 2^p$ (?) [Upfal/STOC 1989, Leighton/Leiserson/Kravets MIT/LCS/RSS 8 1990]⁵⁰

Penn ESE680-002 Spring2007 -- DeHon

Constraints

- Total load should not exceed α of net
 - L =mapping ratio (light loading factor)
 - $\alpha LW =$ number into each subtree
 - $\alpha LW \geq W/2^p$
 - $L \geq 1/(2^p\alpha)$



- Cannot expand past the size of subtree
 - $\beta\alpha N \leq N/2^p$
 - $\beta\alpha \leq 1/2^p$

Penn ESE680-002 Spring2007 -- DeHon

51

Extra Switches

- Extra switch factor: $d \times L$

- Try:

- $\beta=2, \alpha=1/10$
 - $d=8$
 - $dL \approx 40$ ($p=1$)

- Try:

- $\beta=1.01, \alpha=1/4 \rightarrow d=6, L \approx 2, dL \approx 12$ ($p=1$)
- $\beta=1.01, \alpha=1/4 \rightarrow d=6, L \approx 2.8, dL \approx 17$ ($p=0.5$)

Penn ESE680-002 Spring2007 -- DeHon

52

Putting it Together

- Base, linear-population trees have $O(N)$ switches
- Make larger by a factor of L (linear factor)
- Dilated version have a factor of d more switches
- Randomly wired expander
 - Can have $O(N)$ switches
 - Guarantee routes
 - Constants < 100 (looks like < 20)
 - Open: how tight can make it?

Penn ESE680-002 Spring2007 -- DeHon

53

Big Ideas [MSB Ideas]

- In addition to wires, must have switches
 - Have significant area and delay
- Rent's Rule locality reduces
 - both wiring and switching requirements
- Naïve switches match wires at $O(N^{2p})$
 - switch area \gg wire area
 - prevent benefit from multiple layers of metal

Penn ESE680-002 Spring2007 -- DeHon

54

Big Ideas [MSB Ideas]

- Can achieve $O(N)$ switches
 - plausibly $O(N)$ area with sufficient metal layers
- Switchbox depopulation
 - save considerably on area (delay)
 - will waste wires
 - May still come out ahead ([evidence to date](#))