

# ESE534: Computer Organization

Day 19: April 5, 2010  
Interconnect 4: Switching

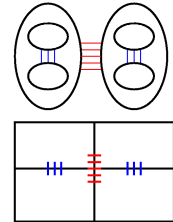
## Previously

- Used Rent's Rule characterization to understand wire growth

$$IO = c N^p$$

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

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



## 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

## Today

- Switching
  - Implications
  - Options
- Exploiting Multiple Metal Layer

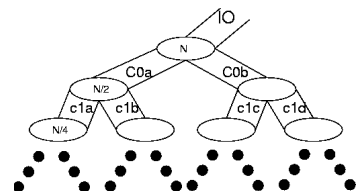
## Switching:

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

## Observation

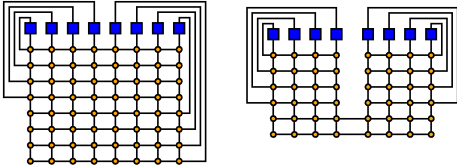
- Locality that saved us wiring, also saves us switching

$$IO = c N^p$$



## Consider

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

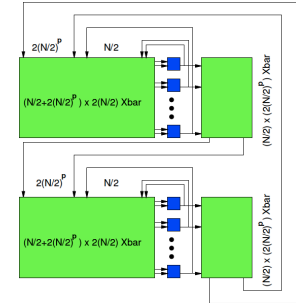


Penn ESE534 Spring 2010 -- DeHon

7

## Preclass 1

- Crosspoints?
- Symbolic ratio?
- Ratio  $N=2^{15}$ ,  $p=2/3$  ?



Penn ESE534 Spring 2010 -- DeHon

## What next

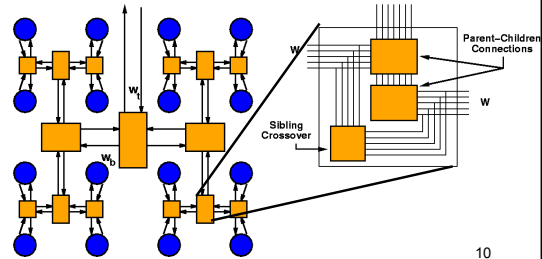
- When something works once?
- ...we try it again...

Penn ESE534 Spring 2010 -- DeHon

9

## Recurse

- Repeat at each level  $\rightarrow$  form a tree

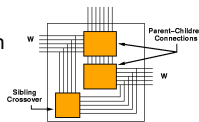


Penn ESE534 Spring 2010 -- DeHon

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  $O(N^{2p})$ 
      - $2 W_{top} \times W_{bot} + (W_{bot})^2$
      - $2 (N^p \times (N/2)^p) + (N/2)^{2p}$
      - ~~$N^{2p}(1/2^p + 1/2^{2p})$~~  Correction from lecture
      - $N^{2p}(2/2^p + 1/2^{2p})$

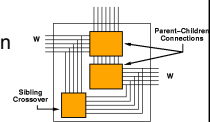


Penn ESE534 Spring 2010 -- DeHon

11

## 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  $O(N^{2p})$ 
      - $N^{2p}(2/2^p + 1/2^{2p})$
    - switches at one level down is
      - $2 \times (N/2)^{2p}(2/2^p + 1/2^{2p})$
      - $2 \times (1/2^p)^2 \times (N^{2p}(2/2^p + 1/2^{2p}))$
      - $2 \times (1/2^p)^2 \times$  previous level

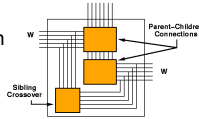


Penn ESE534 Spring 2010 -- DeHon

12

## 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



Now believe this is correct. Just got confused. In lecture.

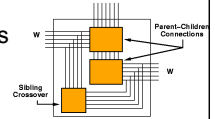
$-2 \times (1/2^p)^2 \times \text{previous level}$   
 $-(2/2^{2p}) = 2^{(1-2p)}$   
 Example:  $p=2/3$   
 $2^{(1-2 \times 2/3)} = 2^{-1/3} = 0.7937$   
 – coefficient < 1 for  $p > 0.5$

13

Penn ESE534 Spring 2010 -- DeHon

## 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 ESE534 Spring 2010 -- DeHon

14

## Good News

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

Penn ESE534 Spring 2010 -- DeHon

15

## Bad News

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

Penn ESE534 Spring 2010 -- DeHon

16

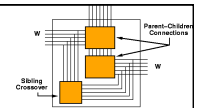
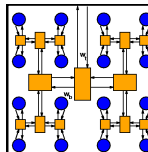
## Additional Structure

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

Penn ESE534 Spring 2010 -- DeHon

17

## Can we do better?



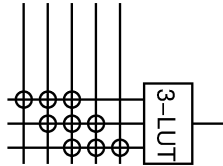
- Crossbar too powerful?
  - Does the specific down channel matter?
- What do we want to do?
  - Connect to any channel on lower level
  - Choose a subset of wires from upper level
    - order not important

Penn ESE534 Spring 2010 -- DeHon

18

## N choose K

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

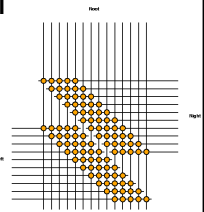


Penn ESE534 Spring 2010 -- DeHon

19

## 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 ESE534 Spring 2010 -- DeHon

20

## Multistage Switching

Penn ESE534 Spring 2010 -- DeHon

21

## Multistage Switching

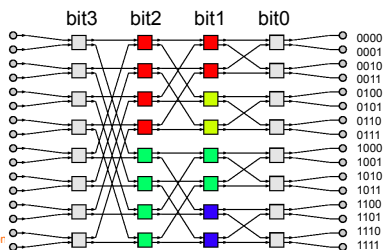
- We can route any **permutation** with fewer switches than a crossbar
- If we allow switching in stages
  - Trade increase in switches in path
  - For decrease in total switches

Penn ESE534 Spring 2010 -- DeHon

22

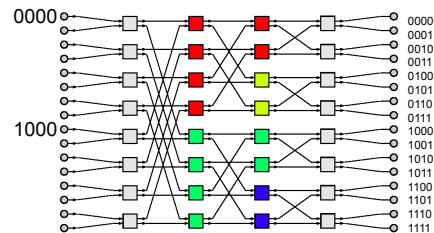
## Butterfly

- Log stages
- Resolve one bit per stage



Penn ESE534 Spring 2010 -- DeHon

## What can a Butterfly Route?

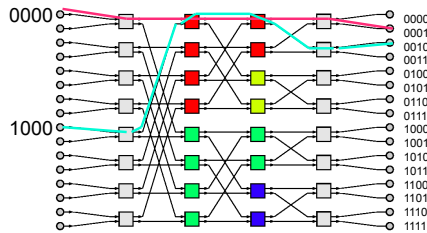


- 0000  $\rightarrow$  0001
- 1000  $\rightarrow$  0010

Penn ESE534 Spring 2010 -- DeHon

24

## What can a Butterfly Route?



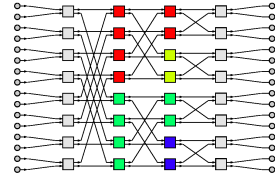
- 0000 → 0001
- 1000 → 0010

Penn ESE534 Spring 2010 -- DeHon

25

## Butterfly Routing

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

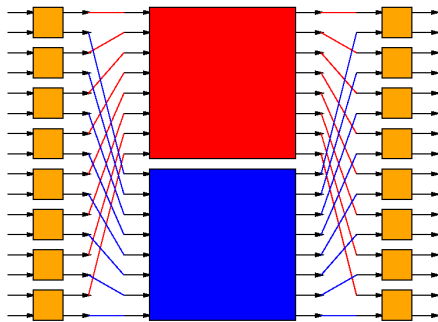


- What required for non-blocking network?

Penn ESE534 Spring 2010 -- DeHon

26

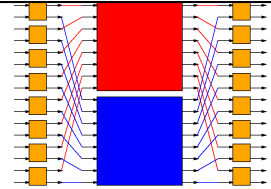
## Decomposition



Penn ESE534 Spring 2010 -- DeHon

27

## Decomposed Routing



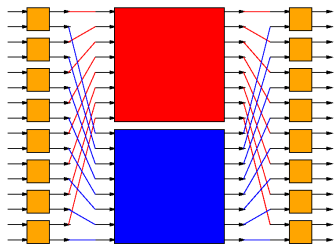
- Pick a link to route.
- Route to destination over **red** network
- At destination,
  - What can we say about the link which shares the final stage switch with this one?
  - What can we do with link?
- Route that link
  - What constraint does this impose?
  - So what do we do?

Penn ESE534 Spring 2010 -- DeHon

28

## Decomposition

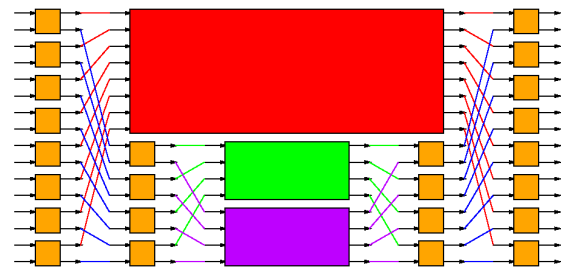
- Switches:  $N/2 \times 2 \times 4 + (N/2)^2 < N^2$



Penn ESE534 Spring 2010 -- DeHon

## Recurse

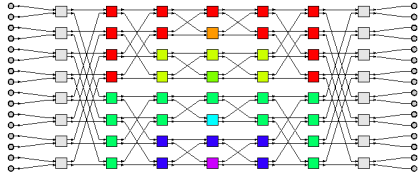
If it works once, try it again...



Penn ESE534 Spring 2010 -- DeHon

## Result: Beneš Network

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

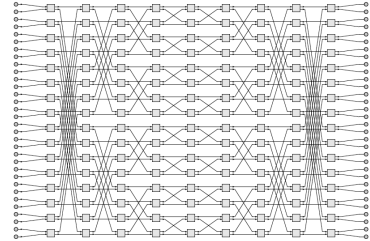


Penn ESE534 Spring 2010 -- DeHon

31

## Preclass 2

- Switches in Beneš 32?
- Ratio to  $32\times 32$  Crossbar?

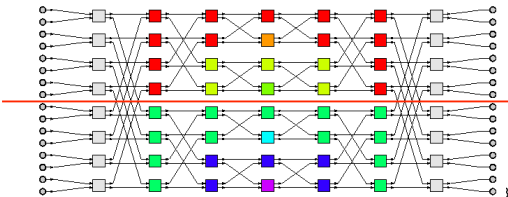


Penn ESE534 Spring 2010 -- DeHon

32

## Beneš Network Wiring

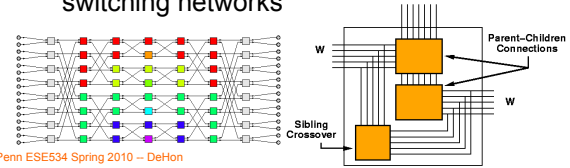
- Bisection:  $N$
- Wiring  $\rightarrow O(N^2)$  area (fixed wire layers)



Penn ESE534 Spring 2010 -- DeHon

## Beneš Switching

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



Penn ESE534 Spring 2010 -- DeHon

## Beneš Switching

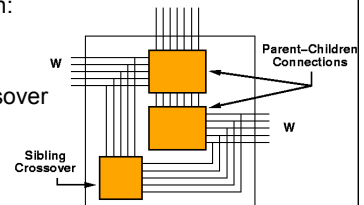
- Implication of Beneš Switching
  - still require  $O(W^2)$  wiring per tree node
    - or a total of  $O(N^{2p})$  wiring
  - now  $O(W \log(W))$  switches per tree node
    - converges to  $O(N)$  total switches!
  - $O(\log^2(N))$  switches in path across network
    - strictly speaking, dominated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

Penn ESE534 Spring 2010 -- DeHon

35

## 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 ESE534 Spring 2010 -- DeHon

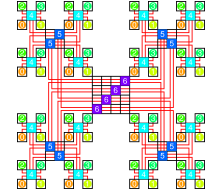
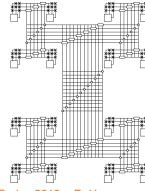
## Linear Switch Population

Penn ESE534 Spring 2010 -- DeHon

37

## Linear Switch Population

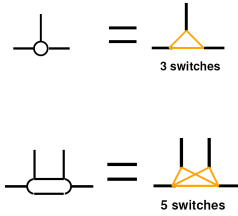
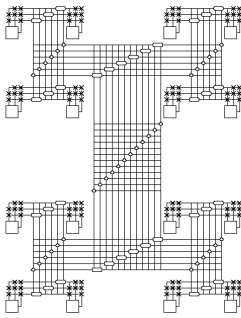
- Can further reduce switches
  - connect each lower channel to  $O(1)$  channels in each tree node
  - end up with  $O(W)$  switches per tree node



Penn ESE534 Spring 2010 -- DeHon

38

## Linear Switch ( $p=0.5$ )

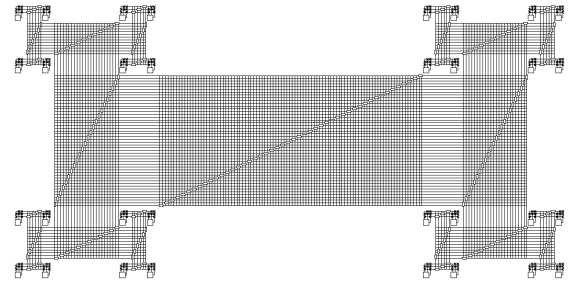


Penn ESE534 Spring 2010 -- DeHon

39

## Linear Population and Beneš

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

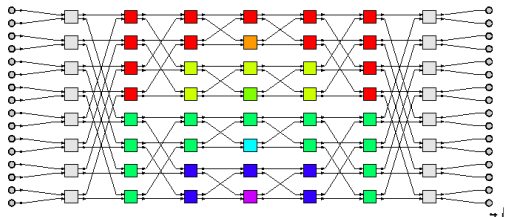


Penn ESE534 Spring 2010 -- DeHon

40

## Beneš Compare

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



Penn ESE534 Spring 2010 -- DeHon

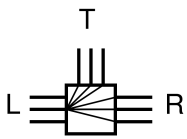
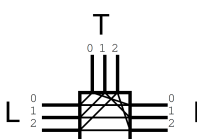
## Linear Consequences: Good News

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

Penn ESE534 Spring 2010 -- DeHon

42

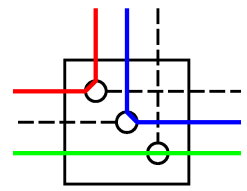
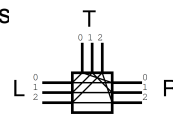
### Preclass 3

- Route Sets:
  - $T \rightarrow (L \& R), L \rightarrow (T \& R), R \rightarrow (T \& L)$
  - $T \rightarrow L, R \rightarrow T, L \rightarrow R$
  - $T \rightarrow L, L \rightarrow R, R \rightarrow L$

Penn ESE534 Spring 2010 -- DeHon 43

### Mapping Ratio

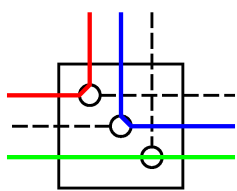
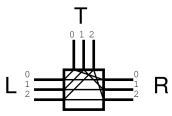



- Mapping ratio says
  - if I have  $W$  channels
    - may only be able to use  $W/MR$  wires
      - for a particular design's connection pattern
    - to accommodate any design
      - for all channels

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

Penn ESE534 Spring 2010 -- DeHon 44

### Mapping Ratio

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

Penn ESE534 Spring 2010 -- DeHon 45

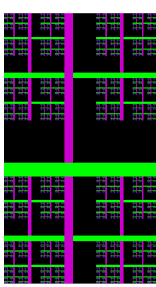

### Linear Consequences: Bad News

- 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 ESE534 Spring 2010 -- DeHon 46

### Area Comparison

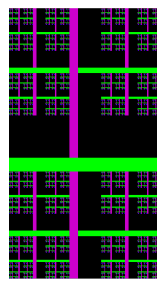

Both:  
 $p=0.67$   
 $N=1024$

M-choose-N  
perfect map
Linear  
MR=2

Penn ESE534 Spring 2010 -- DeHon 47

### Area Comparison

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

M-choose-N  
perfect map
Linear  
MR=2

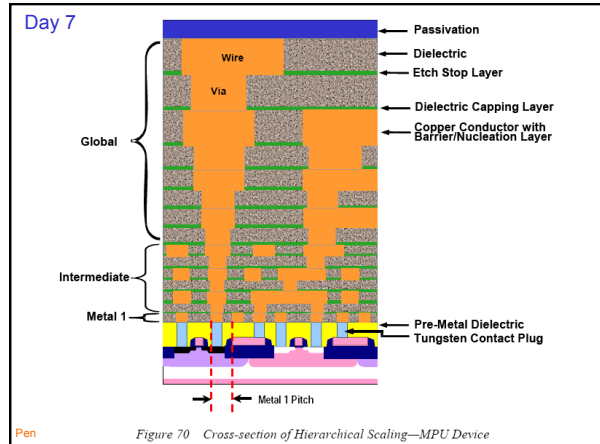
Penn ESE534 Spring 2010 -- DeHon 48



# Multilayer Metal

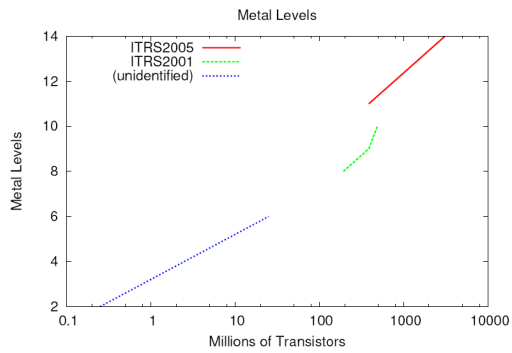
Penn ESE534 Spring 2010 -- DeHon

49



Pen

# Day 7 Wire Layers = More Wiring



Penn ESE534 Spring 2010 -- DeHon

51

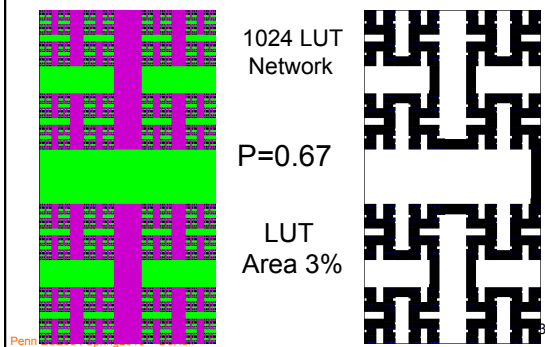
# Opportunity

- Multiple Layers of metal allow us to
  - Increase effective pitch
  - Potentially route in 3D volume

Penn ESE534 Spring 2010 -- DeHon

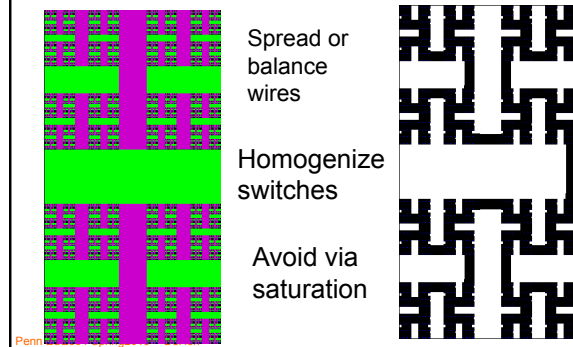
52

# Day 18 Larger “Cartoon”



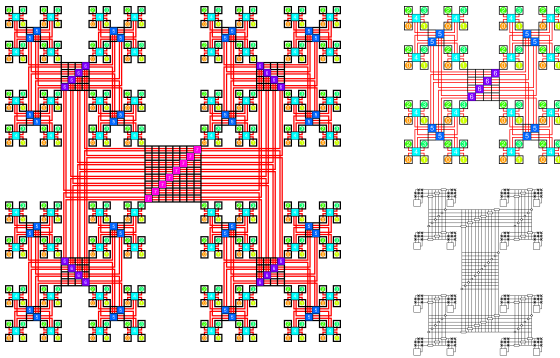
Penn

# Challenge

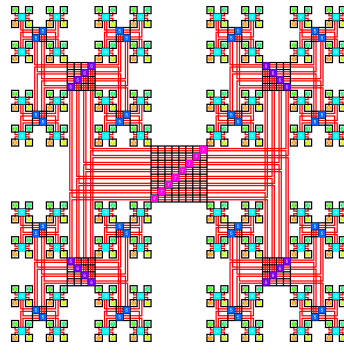


Penn

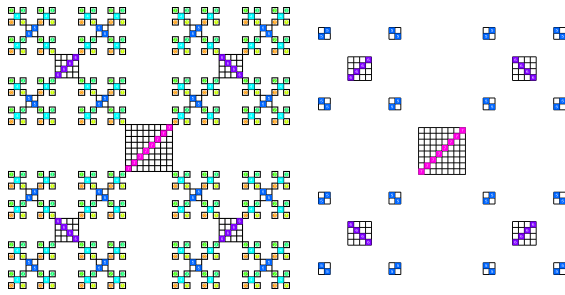
Linear Population Tree (P=0.5)



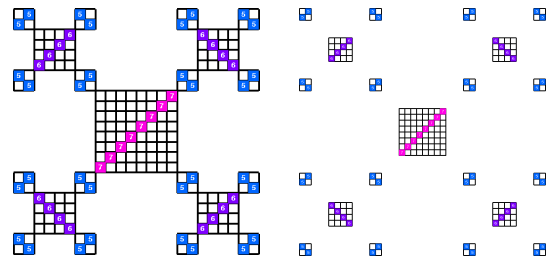
Linear Population Tree (P=0.5)



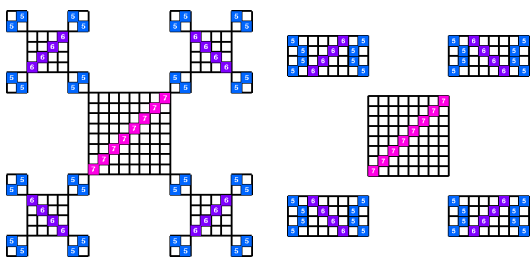
BFT Folding



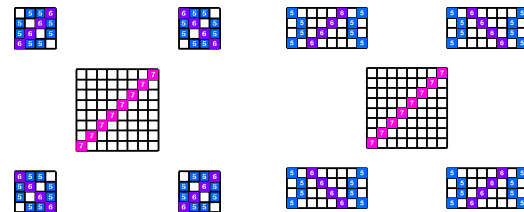
BFT Folding



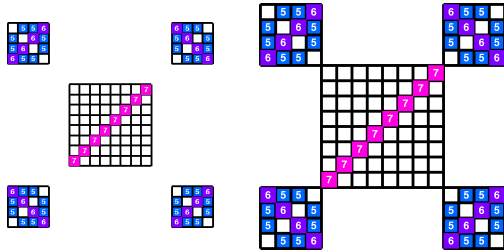
BFT Folding



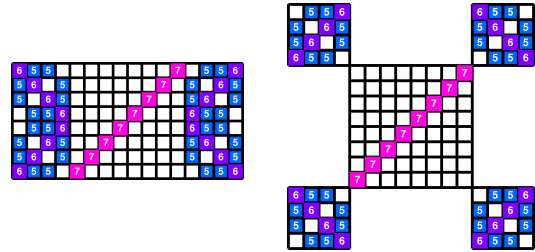
BFT Folding



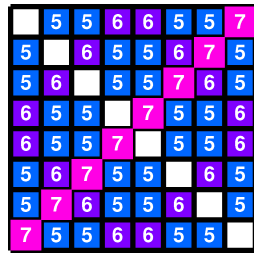
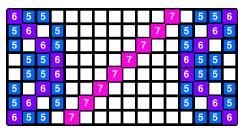
### BFT Folding



### BFT Folding

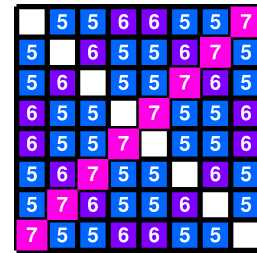


### BFT Folding

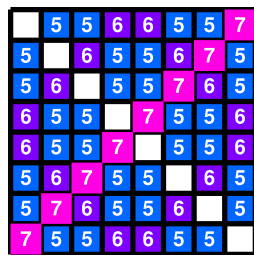
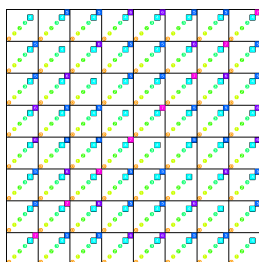


### Invariants

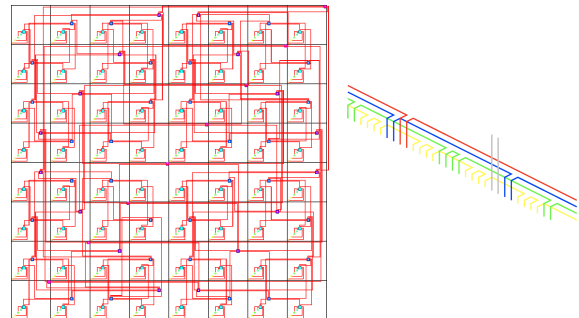
- Lower folds leave both diagonals free
- Current level consumes one, leaving other free



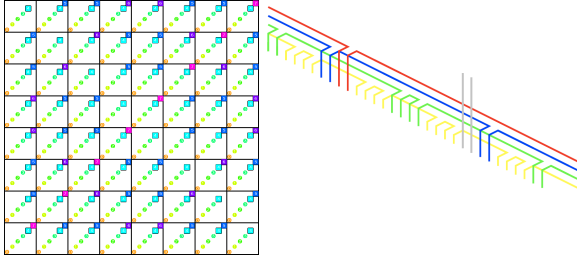
### BFT Folding



### Wiring

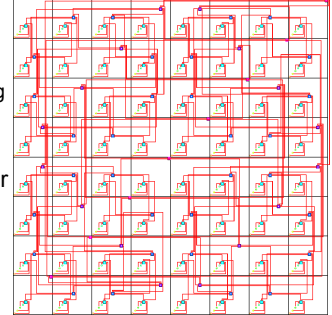


## Avoid via Saturation



## Compact, Multilayer Linear Population Tree Layout

- Can layout BFT
  - in  $O(N)$  2D area
  - with  $O(\log(N))$  wiring layers
- Can be extended for  $p > 0.5$  as well
  - Wire layers grow as  $O(N^{(p-0.5)})$



## Admin

- HW6 Graded
- HW8 dues 4/12
  - Parts still need Wednesday's lecture
- Reading for Wednesday on web

## Big Ideas [MSB Ideas]

- In addition to wires, must have switches
  - 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

## 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)