

ESE534: Computer Organization

Day 17: March 29, 2010
Interconnect 2: Wiring
Requirements and Implications



Previously

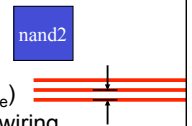
- Identified need for Interconnect
- Seen that interconnect can be expensive
- Identified need to understand/exploit **structure** in our interconnect design

Today

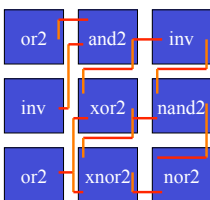
- Wiring Requirements
- Rent's Rule
 - A model of structure
- Implications

Wires and VLSI

- Simple VLSI model
 - Gates have fixed size (A_{gate})
 - Wires have finite spacing (W_{wire})
 - Have a small, finite number of wiring layers
 - *E.g.*
 - one for horizontal wiring
 - one for vertical wiring
 - Assume wires can run over gates



Visually: Wires and VLSI



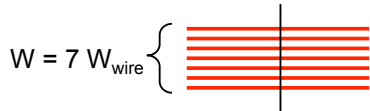
Preclass 1

- How many $40F \times 40F$ gates in $25,000F \times 25,000F$ region?
- How many wires can go in and out?
- Ratio?

Important Consequence

- A set of wires
- crossing a line
- take up space:

$$W = (N \times W_{\text{wire}}) / N_{\text{layers}}$$



Penn ESE534 Spring2010 -- DeHon

7

Thompson's Argument

- The minimum area of a VLSI component is bounded by the larger of:
 - The area to hold all the gates
 - $A_{\text{chip}} \geq N \times A_{\text{gate}}$
 - The area required by the wiring
 - $A_{\text{chip}} \geq N_{\text{horizontal}} W_{\text{wire}} \times N_{\text{vertical}} W_{\text{wire}}$

Penn ESE534 Spring2010 -- DeHon

8

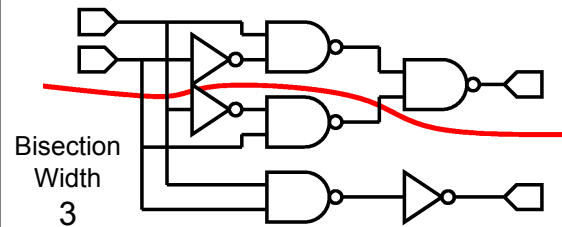
How many wires?

- We can get a **lower bound** on the total number of horizontal (vertical) wires by considering the **bisection** of the computational graph:
 - Cut the graph of gates in half
 - Minimize connections between halves
 - Count number of connections in cut
 - Gives a lower bound on number of wires

Penn ESE534 Spring2010 -- DeHon

9

Bisection

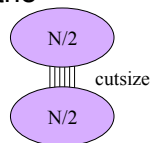


Penn ESE534 Spring2010 -- DeHon

10

Next Question

- In general, if we:
 - Cut design in half
 - Minimizing cut wires
- How many wires will be in the bisection?



Penn ESE534 Spring2010 -- DeHon

11

Arbitrary Graph

- Graph with N nodes
- Cut in half
 - N/2 gates on each side
- **Worst-case?**
 - Every gate output on each side
 - Is used somewhere on other side
 - Cut contains N wires

Penn ESE534 Spring2010 -- DeHon

12

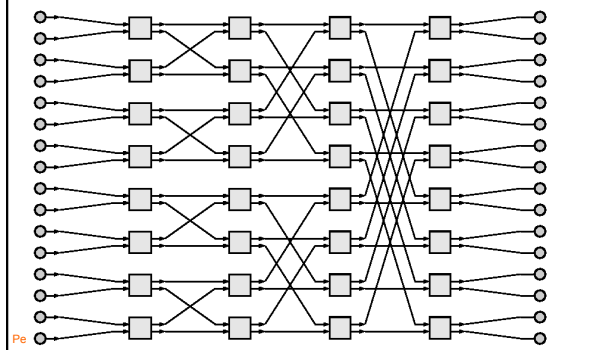
Arbitrary Graph

- For a random graph
 - Something proportional to this is likely
- That is:
 - Given a random graph with N nodes
 - The number of wires in the bisection is likely to be: $c \times N$

Particular Computational Graphs

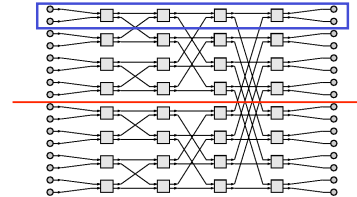
- Some important computations have exactly this property
 - FFT (Fast Fourier Transform)
 - Sorting

FFT



FFT

- Can implement with $N/2$ nodes
 - Group row together
- Any bisection will cut $N/2$ wire bundles
 - True for any reordering



Assembling what we know

- $A_{\text{chip}} \geq N \times A_{\text{gate}}$
- $A_{\text{chip}} \geq N_{\text{horizontal}} W_{\text{wire}} \times N_{\text{vertical}} W_{\text{wire}}$
- $N_{\text{horizontal}} = c \times N$
- $N_{\text{vertical}} = c \times N$
 - [bound true recursively in graph]
- $A_{\text{chip}} \geq cN W_{\text{wire}} \times cN W_{\text{wire}}$

Assembling ...

- $A_{\text{chip}} \geq N \times A_{\text{gate}}$
- $A_{\text{chip}} \geq cN W_{\text{wire}} \times cN W_{\text{wire}}$
- $A_{\text{chip}} \geq (cN W_{\text{wire}})^2$
- $A_{\text{chip}} \geq N^2 \times c'$

Result

- $A_{\text{chip}} \geq N \times A_{\text{gate}}$
- $A_{\text{chip}} \geq N^2 \times c'$
- Wire area grows faster than gate area
- Wire area grows with the square of gate area
- For sufficiently large N,
 - Wire area dominates gate area

Preclass 2

- How does ratio change for 100,000 F × 100,000 F region?

Intuitive Version

- Consider a region of a chip
- Gate capacity in the region goes as area (s^2)
- Wiring capacity into region goes as perimeter ($4s$)
- Perimeter grows more slowly than area
 - Wire capacity saturates before gate



Result

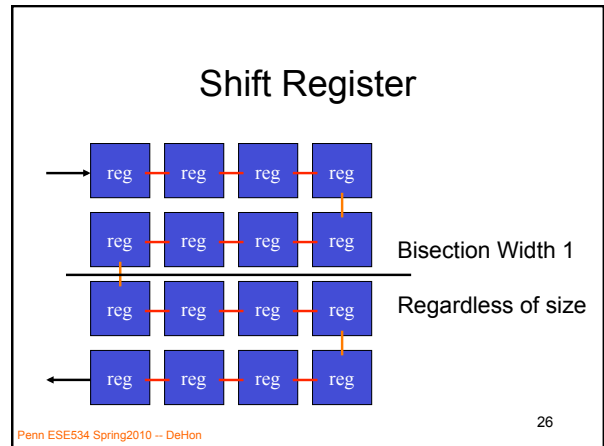
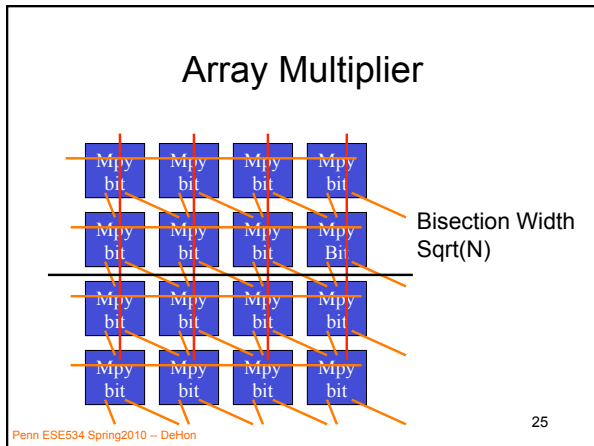
- $A_{\text{chip}} \geq N^2 \times c'$
- Wire area grows with the square of gate area
- Troubling:
 - To **double** the size of our computation
 - Must **quadruple** the size of our chip!

So what?

What do we do with this observation?

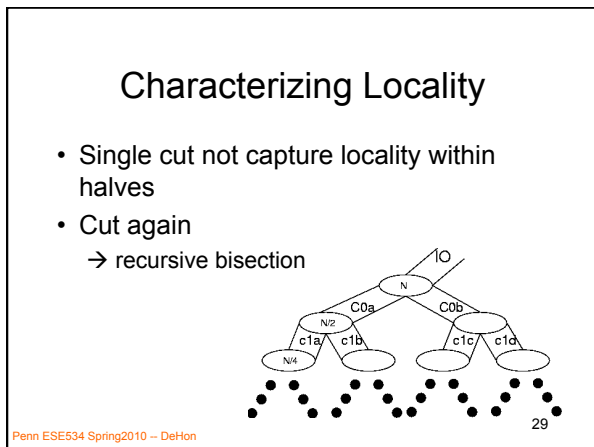
First Observation

- Not all designs have this large of a bisection
- What is typical?



- ### Architecture \Leftrightarrow Structure
- Typical architecture trick:
 - exploit expected problem structure
 - What structure do we have?
 - Impact on resources required?
- Penn ESE534 Spring2010 -- DeHon 27

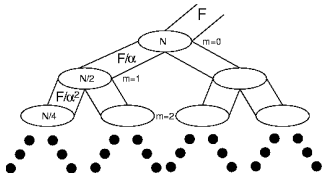
- ### Bisection Bandwidth
- Bisection bandwidth of design
 - lower bound on wire crossings
 - important, **first order** property of a design.
 - Measure to characterize
 - Rather than assume worst case
 - Design with more locality
 - lower bisection bandwidth
 - Enough?
-
- Penn ESE534 Spring2010 -- DeHon 28



- ### Regularizing Growth
- How do bisection bandwidths shrink (grow) at different levels of bisection hierarchy?
 - Basic assumption: Geometric
 - 1
 - $1/\alpha$
 - $1/\alpha^2$
- Penn ESE534 Spring2010 -- DeHon 30

Geometric Growth

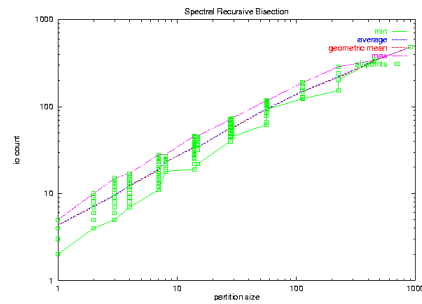
- (F, α) -bifurcator
 - F bandwidth at root
 - geometric regression α at each level



Penn ESE534 Spring2010 -- DeHon

31

Good Model?



Log-log plot → straight lines represent geometric growth

Penn ESE534 Spring2010 -- DeHon

32

Rent's Rule

- In the world of circuit design, an empirical relationship to capture:

$$IO = c N^p$$

- $0 \leq p \leq 1$
- p – characterizes interconnect richness
- Typical: $0.5 \leq p \leq 0.7$
- “High-Speed” Logic $p=0.67$

Penn ESE534 Spring2010 -- DeHon

33

Rent's Rule

- In the world of circuit design, an empirical relationship to capture:

$$IO = c N^p$$

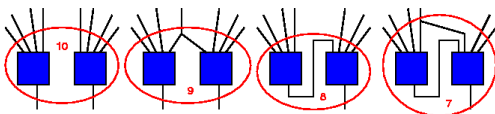
- compare (F, α) -bifurcator
 $\alpha = 2^p$

Penn ESE534 Spring2010 -- DeHon

34

Rent and Locality

- Rent and IO quantifying locality
 - local consumption
 - local fanout

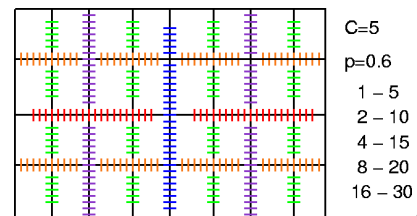


Penn ESE534 Spring2010 -- DeHon

35

What tell us about design?

- Recursive bandwidth requirements in network



Penn ESE534 Spring2010 -- DeHon

36

As a function of Bisection

- $A_{\text{chip}} \geq N \times A_{\text{gate}}$
- $A_{\text{chip}} \geq N_{\text{horizontal}} W_{\text{wire}} \times N_{\text{vertical}} W_{\text{wire}}$
- $N_{\text{horizontal}} = N_{\text{vertical}} = \text{IO} = cN^p$
- $A_{\text{chip}} \geq (cN)^{2p}$
- If $p < 0.5$

$$A_{\text{chip}} \propto N$$

- If $p > 0.5$

$$A_{\text{chip}} \propto N^{2p}$$

Penn ESE534 Spring2010 -- DeHon

37

In terms of Rent's Rule

- If $p < 0.5$, $A_{\text{chip}} \propto N$
- If $p > 0.5$, $A_{\text{chip}} \propto N^{2p}$
- **Typical designs have $p > 0.5$**
→ **interconnect dominates**

Penn ESE534 Spring2010 -- DeHon

38

What tell us about design?

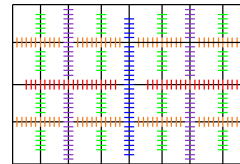
- Recursive bandwidth requirements in network
 - **lower bound** on resource requirements
- N.B. **necessary** but not **sufficient** condition on network design
 - *i.e.* design must also be able to *use* the wires

Penn ESE534 Spring2010 -- DeHon

39

What tell us about design?

- Interconnect lengths
 - Intuition
 - if $p > 0.5$, everything cannot be nearest neighbor
 - as p grows, so wire distances



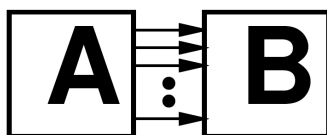
Can think of p as dimensionality:
 $p = 1 - 1/d$

Penn ESE534 Spring2010 -- DeHon

40

Preclass 3

- 25,000 F side, $40F \times 40F$ gates
- Wire length?

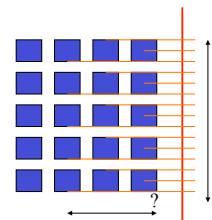


Penn ESE534 Spring2010 -- DeHon

41

Preclass 3

- What's minimum length for longest wires?



Penn ESE534 Spring2010 -- DeHon

42

Generalizing Interconnect Lengths

- $P > 0.5$
- Side is \sqrt{N}
- IO crossing it is N^P
- What's minimum length for longest wires?
- Implication:
 - Wire lengths grow at least as fast as $N^{(p-0.5)}$

$BW = N^P$

\sqrt{N}

$N^{(p-0.5)}$

Penn ESE534 Spring2010 -- DeHon 43

Delays

Recall from Day 7

- Logical capacities growing
- Wirelengths?
 - No locality $\propto \kappa$
 - Rent's Rule
 - $L \propto n^{(p-0.5)}$
 - $[p > 0.5]$

Relative Logical Capacity

Penn ESE534 Spring2010 -- DeHon

Capacity

- Rent: $IO = C * N^P$
- $p > 0.5$
- $A = C * N^{2p}$
- $N = (A/C)^{(1/2p)}$
- Logical Area $\propto \kappa^2$
- $N' = ((\kappa^2 A)/C)^{(1/2p)}$
- $N' = (A/C)^{(1/2p)} \times (\kappa^2)^{(1/2p)}$
- $N' = N \times (\kappa^2)^{(1/2p)}$
- $N' = N \times (\kappa)^{(1/p)}$

- Sanity Check
 - $p = 1$
 - $N_2 = \kappa N$
 - $p \sim 0.5$
 - $N_2 \sim \kappa^2 N$

Penn ESE534 Spring2010 -- DeHon 45

What tell us about design?

- $IO \propto N^P$
- Bisection $BW \propto N^P$
- side length $\propto N^P$
 - N if $p < 0.5$
- Area $\propto N^{2p}$
 - $p > 0.5$
- Average Wire Length $\propto N^{(p-0.5)}$
 - $p > 0.5$

N.B. 2D VLSI world has "natural" Rent of $P = 0.5$ (area vs. perimeter)

Penn ESE534 Spring2010 -- DeHon 46

Preclass 4

- Depth 20 circuit, 2-input gates
 - Maximum number of gates?
 - Topology?
 - Rent p ?
 - Minimum distance?
 - Lower bound maximum length
- Depth 24 circuit
 - Lower bound maximum length?

Penn ESE534 Spring2010 -- DeHon 47

Rent's Rule Caveats

- Modern "systems" on a chip -- likely to contain subcomponents of varying Rent complexity
- Less I/O at certain "natural" boundaries
- System close
 - Rent's Rule apply to workstation, PC, MP3 player, Smart Phone?

Penn ESE534 Spring2010 -- DeHon 48

Area/Wire Length

- Bad news
 - Area $\sim \Omega(N^{2p})$
 - faster than N
 - Avg. Wire Length $\sim \Omega(N^{(p-0.5)})$
 - grows with N
- Can designers/CAD control p (locality) once appreciate its effects?
- *I.e.* maybe this cost changes design style/criteria so we mitigate effects?

Penn ESE534 Spring2010 -- DeHon

49

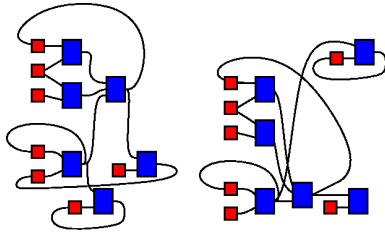
What Rent didn't tell us

- Bisection bandwidth purely geometrical
- No constraint for delay
 - *I.e.* a partition may leave critical path weaving between halves

Penn ESE534 Spring2010 -- DeHon

50

Critical Path and Bisection



Minimum cut may cross critical path multiple times.
Minimizing long wires in critical path \rightarrow increase cut size.

Penn ESE534 Spring2010 -- DeHon

51

Original Memo

- Current Issue (Winter 2010, v2n1) of IEEE Solid-State Circuits Magazine
- Retrospect on IBM 1401 and E. F. Rent
 - Including original memos
- Added link to reading

Penn ESE534 Spring2010 -- DeHon

52

Admin

- HW5 graded
- HW8 out – due April 12th
- Reading for Wed. on web

Penn ESE534 Spring2010 -- DeHon

53

Big Ideas [MSB Ideas]

- Rent's rule characterizes locality
Fixed wire layers:
 - \rightarrow Area growth $\Omega(N^{2p})$
 - \rightarrow Wire Length $\Omega(N^{(p-0.5)})$
- $p > 0.5 \rightarrow$ interconnect growing faster than compute elements
 - expect **interconnect to dominate** other resources

Penn ESE534 Spring2010 -- DeHon

54