

ESE534: Computer Organization

Day 15: March 14, 2012
Interconnect 2: Wiring
Requirements and Implications



Previously

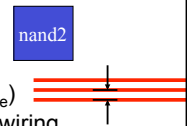
- Identified need for Interconnect
- Seen that interconnect can be expensive
- Identified need to understand/exploit **structure** in our interconnect design

Today

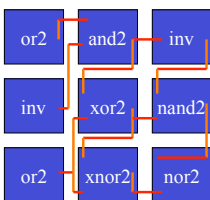
- Wiring Requirements
- Rent's Rule
 - A model of structure
- Implications

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



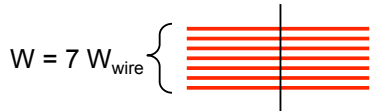
Preclass 1

- How many 40F×40F gates in 25,000F×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 Spring2012 -- 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 Spring2012 -- DeHon

8

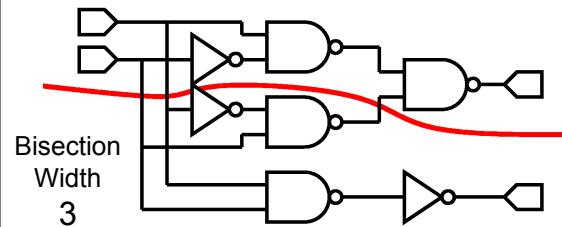
How many wires?

- 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 Spring2012 -- DeHon

9

Bisection

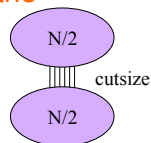


Penn ESE534 Spring2012 -- DeHon

10

Next Question

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



Penn ESE534 Spring2012 -- 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 Spring2012 -- DeHon

12

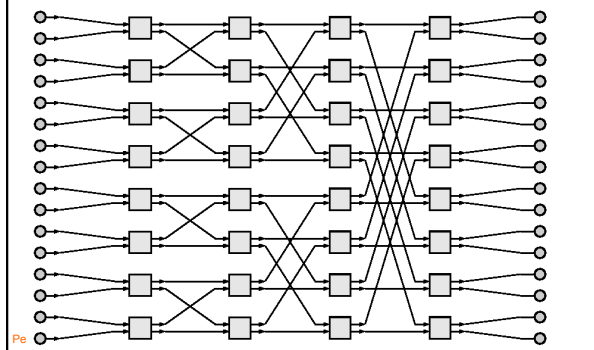
Arbitrary Graph

- 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

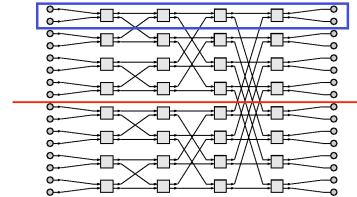
- Some important computations have exactly this property
 - FFT (Fast Fourier Transform)
 - Sorting

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?

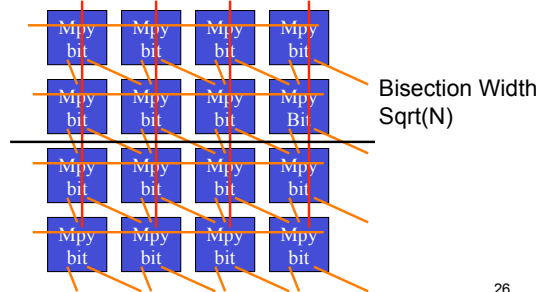
First Observation

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

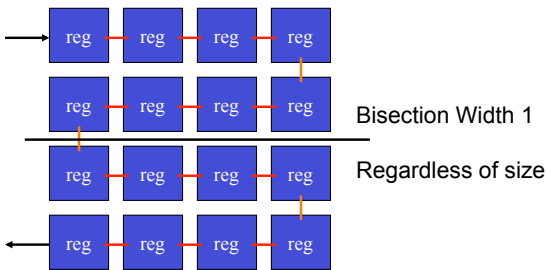
Favorite Design Elements

- What are your favorite computing design elements?
- What are the bisection bandwidths for these elements?

Array Multiplier



Shift Register

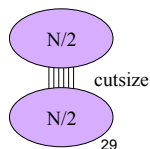


Architecture \leftrightarrow Structure

- Typical architecture trick:
 - exploit expected problem structure
- What structure do we have?
- Impact on resources required?

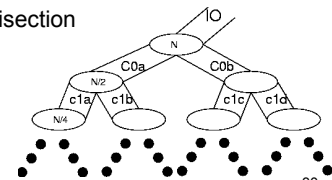
Bisection Bandwidth

- 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?



Characterizing Locality

- Single cut does not capture locality within halves
- Cut again
 - recursive bisection

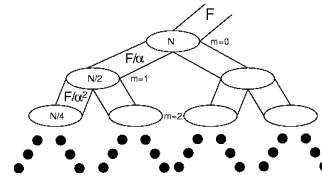


Regularizing Growth

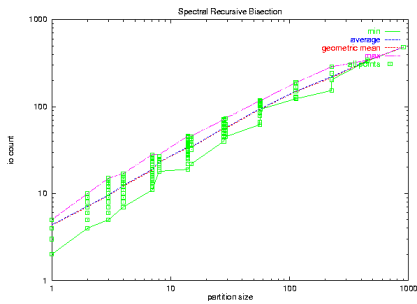
- How do bisection bandwidths shrink (grow) at different levels of bisection hierarchy?
- Basic assumption: Geometric
 - 1
 - $1/\alpha$
 - $1/\alpha^2$

Geometric Growth

- F bandwidth at root
- geometric regression α at each level
 - Or growth by α for every doubling



Good Model?



Log-log plot → straight lines represent geometric growth

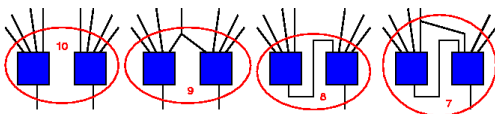
Rent's Rule

- 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 $p=0.67$

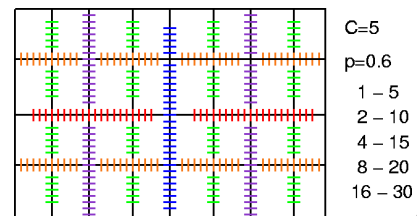
Rent and Locality

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



What tell us about design?

- Recursive bandwidth requirements in network



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 Spring2012 -- 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 Spring2012 -- 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 Spring2012 -- DeHon

39

Capacity Impact

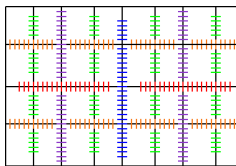
- Rent: $\text{IO} = C \cdot N^p$
- $p > 0.5$
- $A = C \cdot N^{2p}$
- $N = (A/C)^{(1/2p)}$
- Logical Area $\propto (1/S)^2$
- $N' = (((1/S)^2 A)/C)^{(1/2p)}$
- $N' = (A/C)^{(1/2p)} \times ((1/S)^2)^{(1/2p)}$
- $N' = N \times ((1/S)^2)^{(1/2p)}$
- $N' = N \times (1/S)^{(1/p)}$
- Sanity Check
 - $p=1$
 - $N_2 = N/S$
 - $p=0.5$
 - $N_2 = N/S^2$

Penn ESE534 Spring2012 -- DeHon

40

What tell us about design?

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



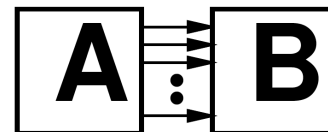
Can think of p as dimensionality:
 $p = 1 - 1/d$

Penn ESE534 Spring2012 -- DeHon

41

Preclass 5

- 24,000 F side, 40F × 40 F gates
- Wire length?

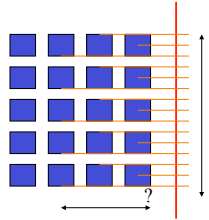


Penn ESE534 Spring2012 -- DeHon

42

Preclass 5

- What's minimum length for longest wires?



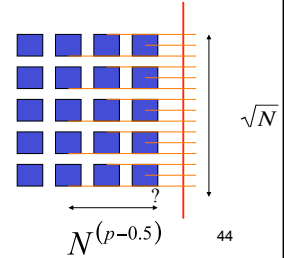
Penn ESE534 Spring2012 -- DeHon

43

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 $N^{(p-0.5)}$

$$BW = N^P$$

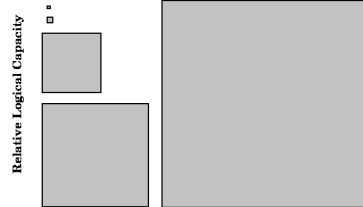


Penn ESE534 Spring2012 -- DeHon

44

Scaling → Delays

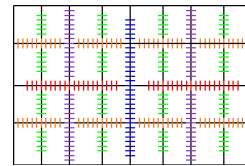
- Logical capacities on chip growing
- Wirelengths?
 - No locality \propto chip-side = $1/S$
 - Rent's Rule
 - $L \propto n^{(p-0.5)}$
 - $[p > 0.5]$



Penn ESE534 Spring2012 -- DeHon

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 Spring2012 -- DeHon

46

Preclass 6

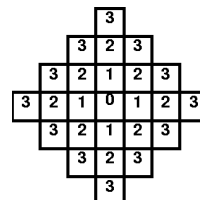
- How many gates reachable with 800F of wiring?
- How many gates reachable with 1600F of wiring?

Penn ESE534 Spring2012 -- DeHon

47

Distance

- How many things at a given distance?



Penn ESE535 Spring 2011 -- DeHon

48

Preclass 7

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

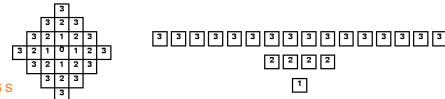
Penn ESE534 Spring2012 -- DeHon

49

“Closeness”

- Try placing “everything” close

Manhattan Distance	Places	Transitive Fanin
1	4	4
2	8	16
3	12	64
i	i	i
n	$4n$	4^n



Penn ESE535 S

50

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 Spring2012 -- DeHon

51

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 Spring2012 -- DeHon

52

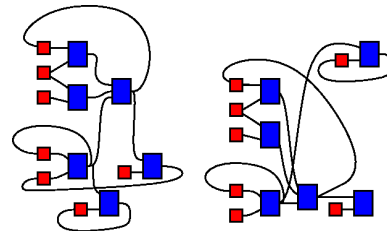
What Rent didn't tell us

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

Penn ESE534 Spring2012 -- DeHon

53

Critical Path and Bisection



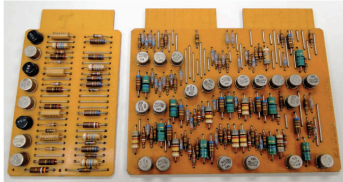
Minimum cut may cross critical path multiple times. Minimizing long wires in critical path \rightarrow increase cut size.

Penn ESE534 Spring2012 -- DeHon

54

Original Memo

- Recent Issue (Winter 2010, v2n1) of IEEE Solid-State Circuits Magazine
- Retrospect on IBM 1401 and E. F. Rent
 - Including original memos
- Linked Supplemental Reading



Penn ESE534 Spring2012 -- DeHon

FIGURE 5: Single- and double-width SMS cards from the IBM 1401 Processing Unit. (Photo courtesy of Robert Garnea)

Admin

- HW7 due Monday
- Reading for Monday on web

Penn ESE534 Spring2012 -- DeHon

56

Big Ideas [MSB Ideas]

- Rent's rule characterizes locality
 - Fixed wire layers:
 - Area growth $\Omega(N^{2p})$
 - 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 Spring2012 -- DeHon

57