

CS184a: Computer Architecture (Structure and Organization)

Day 12: February 5, 2003
Interconnect 2: Wiring
Requirements and Implications



Caltech CS184 Winter2003 -- DeHon

Previously

- Need for Interconnect
- Why simplest things don't work
 - Bus
 - Crossbar
- Need to understand/exploit structure in our interconnect problem

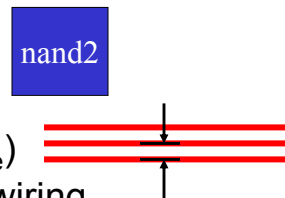
Caltech CS184 Winter2003 -- DeHon

Today

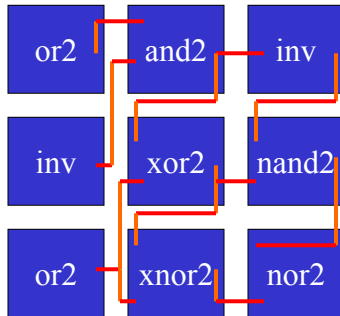
- Wiring Requirements
- Rent's Rule
 - A model of structure
- Implications

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



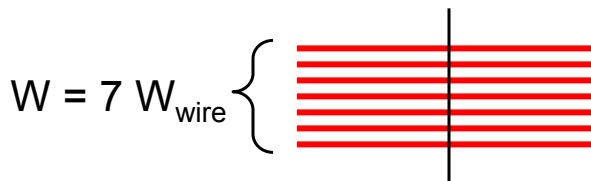
Visually: Wires and VLSI



Important Consequence

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

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



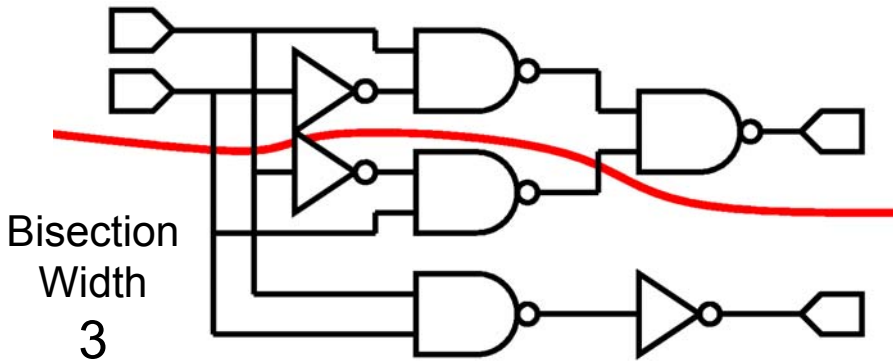
Thompson's Argument

- 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
 - $A_{\text{chip}} \geq N_{\text{horizontal}} W_{\text{wire}} \times N_{\text{vertical}} W_{\text{wire}}$

How many wires?

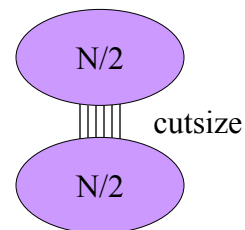
- We can get a **lower bound** on the total number of horizontal (vertical) wires by considering the **bisection** of the computational graph:
 - Cut the graph of gates in half
 - Minimize connections between halves
 - Count number of connections in cut
 - Gives a lower bound on number of wires

Bisection



Next Question

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



Arbitrary Graph

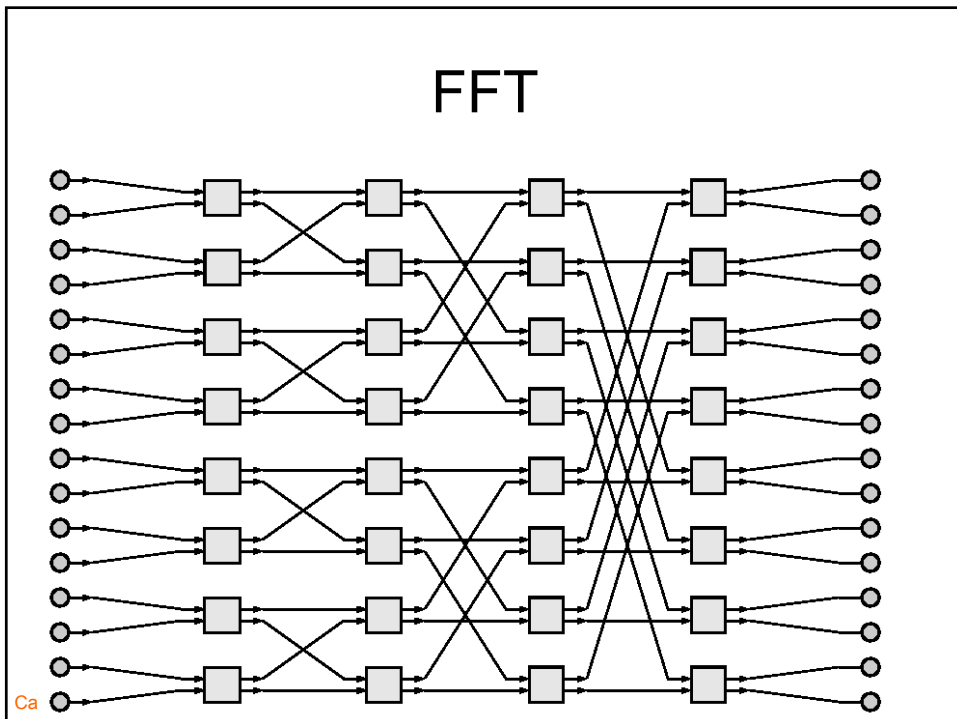
- Graph with N nodes
- Cut in half
 - $N/2$ gates on each side
- **Worst-case:**
 - Every gate output on each side
 - Is used somewhere on other side
 - Cut contains N wires

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: $c \times N$

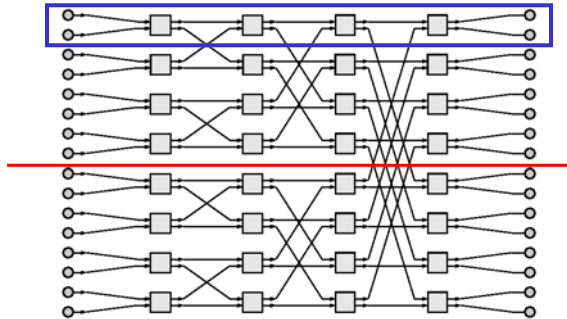
Particular Computational Graphs

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



FFT

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



Caltech CS184 Winter2003 -- DeHon

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]
- $A_{\text{chip}} \geq cN W_{\text{wire}} \times cN W_{\text{wire}}$

Caltech CS184 Winter2003 -- DeHon

Assembling ...

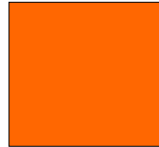
- $A_{\text{chip}} \geq N \times A_{\text{gate}}$
- $A_{\text{chip}} \geq cN W_{\text{wire}} \times cN W_{\text{wire}}$
- $A_{\text{chip}} \geq (cN W_{\text{wire}})^2$
- $A_{\text{chip}} \geq N^2 \times c'$

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

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!

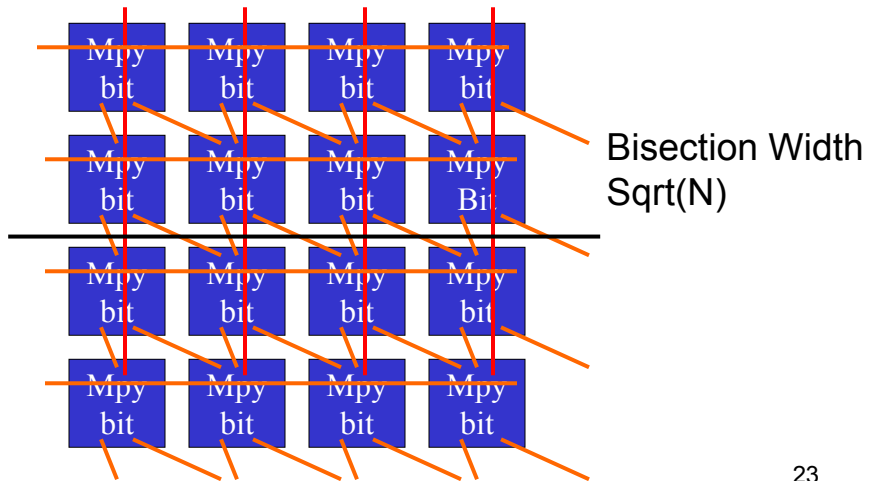
So what?

What do we do with this observation?

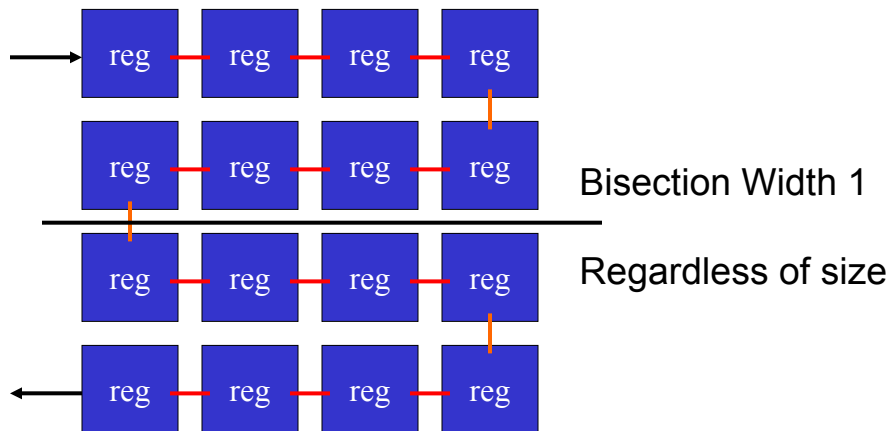
First Observation

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

Array Multiplier



Shift Register

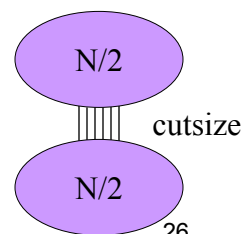


Architecture \Leftrightarrow Structure

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

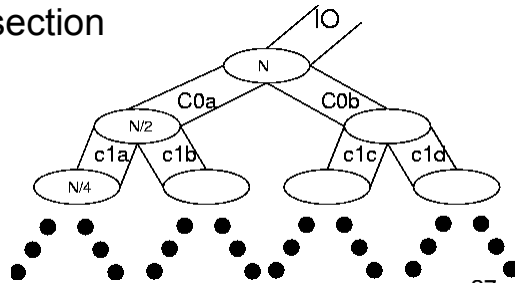
Bisection Bandwidth

- Bisection bandwidth of design
 - lower bound on network bisection bandwidth
 - 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 not capture locality within halves
- Cut again
 - recursive bisection

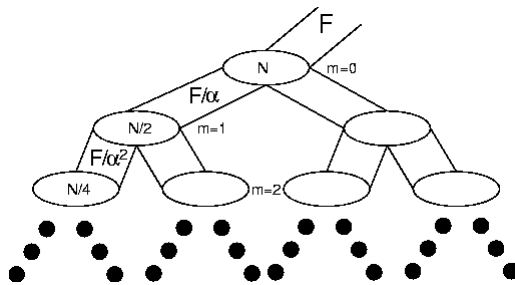


Regularizing Growth

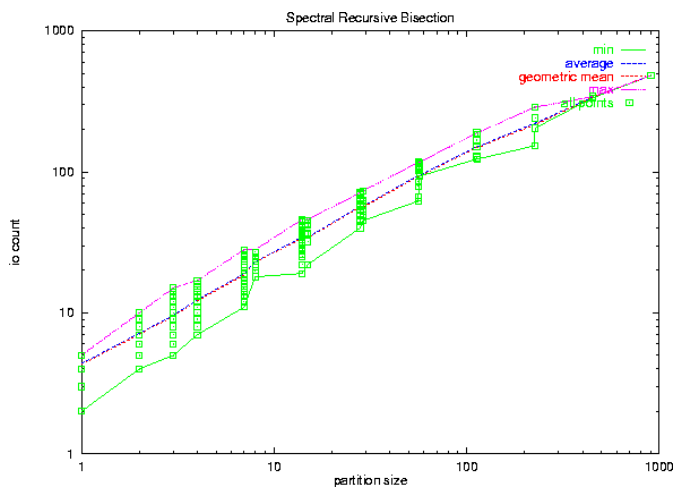
- How do bisection bandwidths shrink (grow) at different levels of bisection hierarchy?
- Basic assumption: Geometric
 - 1
 - $1/\alpha$
 - $1/\alpha^2$

Geometric Growth

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



Good Model?



Log-log plot \rightarrow 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's Rule

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

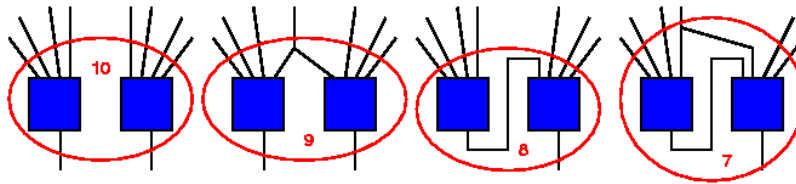
$$IO = c N^p$$

- compare (F, α) -bifurcator

$$\alpha = 2^p$$

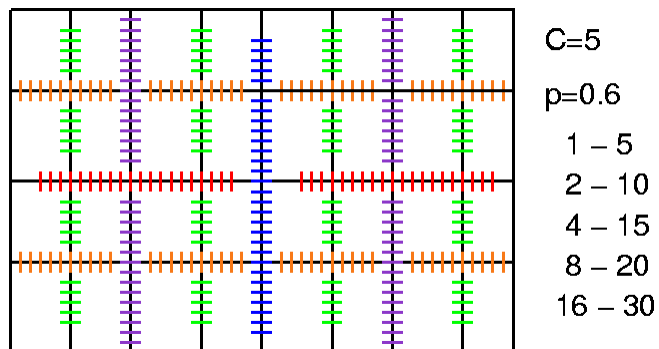
Rent and Locality

- Rent and IO capture/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 $p > 0.5$

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

In terms of Rent's Rule

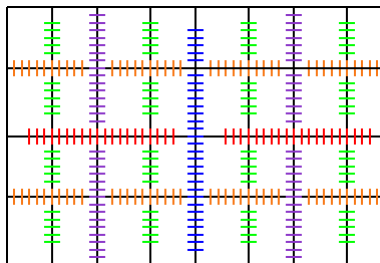
- If $p < 0.5$, $A_{\text{chip}} \propto N$
- If $p > 0.5$, $A_{\text{chip}} \propto N^{2p}$
- **Typical** designs have **$p > 0.5$**
→ **interconnect dominates**

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

What tell us about design?

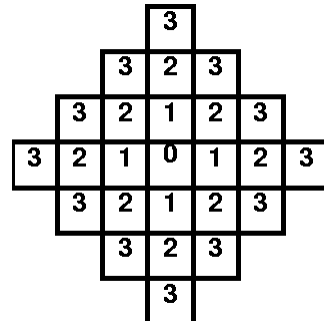
- Interconnect lengths
 - Intuition
 - if $p > 0.5$, everything cannot be nearest neighbor
 - as p grows, so wire distances



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

What tell us about design?

- Interconnect lengths
 - $IO=(n^2)^P$ cross distance n
 - $D(IO)/dn$ end at exactly distance n
 - $E(l)=\text{Integral } 0 \text{ to } n=\sqrt{N}$
 - of $n*(d(IO)/dn)/n^2$
 - assume iid sources



Caltech CS184 Winter2003 -- DeHon

Math

$$\frac{d(IO)}{dn} = \frac{d(cn^{2p})}{dn} = 2pcn^{2p-1}$$

$$\int_0^{\sqrt{N}} \left(\frac{n \times \frac{d(IO)}{dn}}{n^2} \right) dn$$

$$\int_0^{\sqrt{N}} \left(\frac{n \times 2pcn^{2p-1}}{n^2} \right) dn$$

Caltech CS184 Winter2003 -- DeHon

Math continued

$$\int_0^{\sqrt{N}} \left(\frac{n \times 2pcn^{2p-1}}{n^2} \right) dn$$

$$\int_0^{\sqrt{N}} \left(\frac{2pcn^{2p-1}}{n} \right) dn$$

$$\int_0^{\sqrt{N}} (2pcn^{2p-2}) dn$$

Math continued

$$\int_0^{\sqrt{N}} (2pcn^{2p-2}) dn$$

$$\left(\frac{2pcn^{2p-1}}{2p-1} \right) \Big|_0^{\sqrt{N}}$$

Math Continued

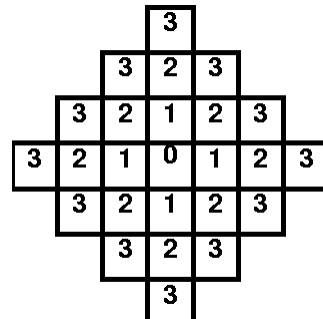
$$\left(\frac{2pcn^{2p-1}}{2p-1} \right) \Big|_0^{\sqrt{N}}$$

$$\left(\frac{2pc}{2p-1} \right) (N^{0.5})^{2p-1}$$

$$\left(\frac{2pc}{2p-1} \right) (N^{p-0.5})$$

What tell us about design?

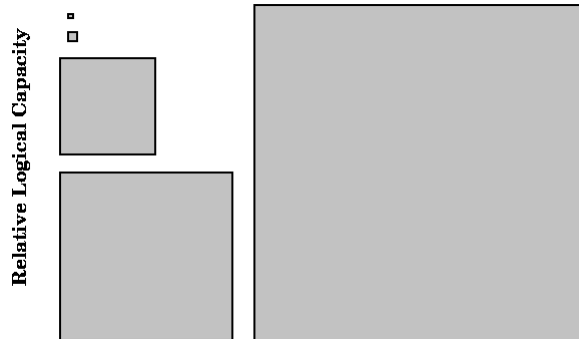
- Interconnect lengths
 - $IO=(n^2)^P$ cross distance n
 - $D(IO)/dn$ end at exactly distance n
 - $E(I)=\text{Integral } 0 \text{ to } n=\sqrt{N}$
 - of $n*(d(IO)/dn)/n^2$
 - assume iid sources
 - $E(I)=\Omega(N^{(p-0.5)})$
 - $p>0.5$



True even with multiple metal layers.

Delays

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



Caltech CS184 Winter2003 -- DeHon

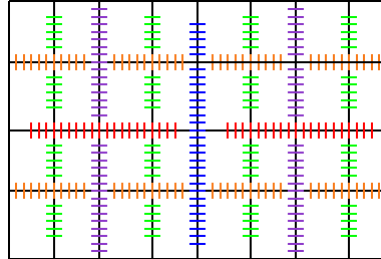
Capacity

- Rent: $IO = C * N^p$
- $p > 0.5$
- $A = C * N^{2p}$
- Logical Area $\rightarrow \kappa^2$
 - $\kappa^2 A = C * N_2^{2p}$
 - $\kappa^2 C * N_1^{2p} = C * N_2^{2p}$
 - $\kappa^2 N_1^{2p} = N_2^{2p}$
 - $\kappa N_1^p = N_2^p$
 - $N_2 = \kappa^{(1/p)} N_1$
- Sanity Check
 - $p = 1$
 - $N_2 = \kappa N$
 - $p \sim 0.5$
 - $N_2 \sim \kappa^2 N$

Caltech CS184 Winter2003 -- 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)

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, PDA?)

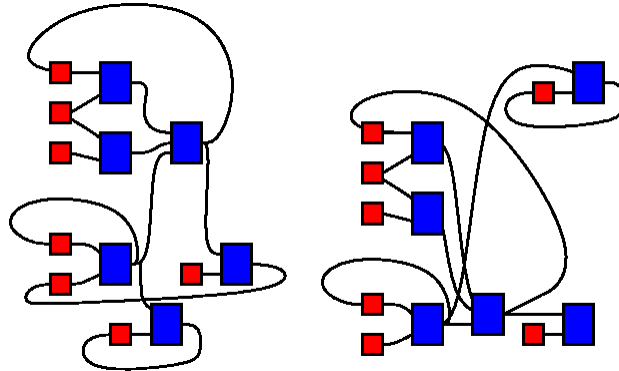
Area/Wire Length

- 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?

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

Critical Path and Bisection



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

Rent Weakness

- Not account for path topology
- ? Can we define a “Temporal” Rent which takes into consideration?
 - Promising research topic

Big Ideas

[MSB Ideas]

- Rent's rule characterize 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