

# CS184b: Computer Architecture (Abstractions and Optimizations)

Day 4: April 4, 2005  
Interconnect



## Previously

- CS184a
  - interconnect needs and requirements
  - basic topology
  - Mostly thought about static/offline routing

## This Quarter

- This quarter
  - parallel systems require
  - typically dynamic switching
  - interfacing issues
    - model, hardware, software

## Today

- Issues
- Topology/locality/scaling
  - (some review)
- Styles
  - from static
  - to online, packet, wormhole
- Online routing

## Issues

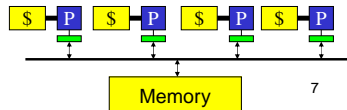
- |   |  |
|---|--|
| <b>Old</b>  | <b>New</b>   |
| <ul style="list-style-type: none"><li>• Bandwidth<ul style="list-style-type: none"><li>– aggregate, per endpoint</li><li>– local contention and hotspots</li></ul></li><li>• Latency</li><li>• Cost (scaling)<ul style="list-style-type: none"><li>– locality</li></ul></li></ul> | <ul style="list-style-type: none"><li>• Arbitration<ul style="list-style-type: none"><li>– conflict resolution</li><li>– deadlock</li></ul></li><li>• Routing<ul style="list-style-type: none"><li>– (quality vs. complexity)</li></ul></li><li>• Ordering (of messages)</li></ul> |

## Topology and Locality

(Partially) Review

## Simple Topologies: Bus

- Single Bus
  - simple, cheap
  - low bandwidth
    - not scale with PEs
  - typically online arbitration
    - can be offline scheduled



Caltech CS184 Spring2005 -- DeHon

7

## Bus Routing

- Offline:
  - divide time into N slots
  - assign positions to various communications
  - run modulo N w/ each consumer/producer send/receiving on time slot
- e.g.
  - 1: A->B
  - 2: C->D
  - 3: A->C
  - 4: A->B
  - 5: C->B
  - 6: D->A
  - 7: D->B
  - 8: A->D

Caltech CS184 Spring2005 -- DeHon

8

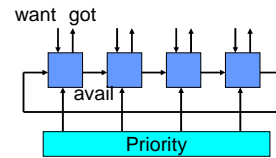
## Bus Routing

- Online:
  - request bus
  - wait for acknowledge
- Priority based:
  - give to highest priority which requests
  - consider ordering
  - $Got_i = Want_i \wedge Avail_i$
  - $Avail_{i+1} = Avail_i \wedge \neg Want_i$
- Solve arbitration in log time using parallel prefix
- For fairness
  - start priority at different node
  - use cyclic parallel prefix
    - deal with variable starting point

Caltech CS184 Spring2005 -- DeHon

9

## Arbitration Logic

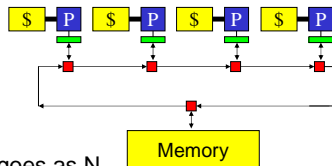


Caltech CS184 Spring2005 -- DeHon

10

## Token Ring

- On bus
  - delay of cycle goes as N
  - can't avoid, even if talking to nearest neighbor
- Token ring
  - pipeline bus data transit (ring)
    - high frequency
  - can exit early if local
  - use token to arbitrate use of bus

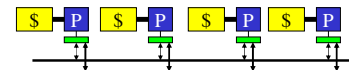


Caltech CS184 Spring2005 -- DeHon

11

## Multiple Busses

- Simple way to increase bandwidth
  - use more than one bus
- Can be static or dynamic assignment to busses
  - static
    - A->B always uses bus 0
    - C->D always uses bus 1
  - dynamic
    - arbitrate for a bus, like instruction dispatch to k identical CPU resources

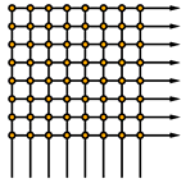


Caltech CS184 Spring2005 -- DeHon

12

## Crossbar

- No bandwidth reduction
  - (except receiver at endpoint)
- Easy routing (on or offline)
- Scales poorly
  - $N^2$  area and delay
- No locality



Caltech CS184 Spring2005 -- DeHon

## Hypercube

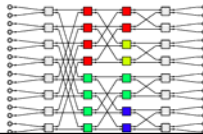
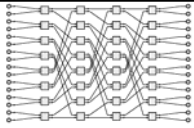
- Arrange  $N=2^n$  nodes in n-dimensional cube
- At most n hops from source to sink
  - $N = \log_2(N)$
- High bisection bandwidth
  - good for traffic (but can you use it?)
  - bad for cost [ $O(n^2)$ ]
- Exploit locality
- Node size grows
  - as  $\log(N)$  [IO]
  - Maybe  $\log^2(N)$  [xbar between dimensions]

Caltech CS184 Spring2005 -- DeHon

14

## Multistage

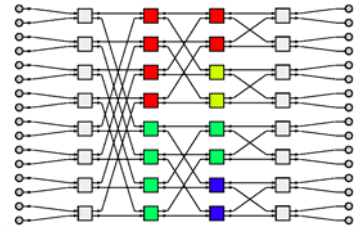
- Unroll hypercube vertices so  $\log(N)$ , constant size switches per hypercube node
  - solve node growth problem
  - lose locality
  - similar good/bad points for rest



Caltech CS184 Spring2005 -- DeHon

## Hypercube/Multistage Blocking

- Minimum length multistage
  - many patterns cause bottlenecks
  - e.g.

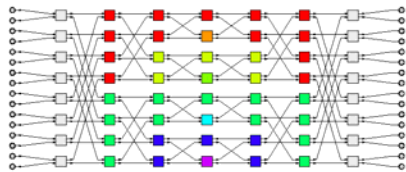


Caltech CS184 Spring2005 -- DeHon

CS184a: Day16

## Beneš Network

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



Caltech CS184 Spring2005 -- DeHon

17

## Online Hypercube Blocking

- If routing offline, can calculate Benes-like route
- Online, don't have time, or global view
- **Observation:** only a few, canonically bad patterns
- **Solution:** Route to random intermediate
  - then route from there to destination
  - ...turns worst-case into average case
    - at the expense of locality

Caltech CS184 Spring2005 -- DeHon

18

## K-ary N-cube

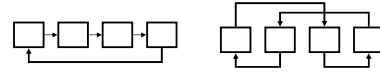
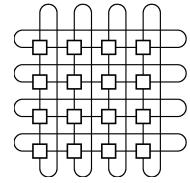
- Alternate reduction from hypercube
  - restrict to  $N < \log(\text{Nodes})$  dimensional structure
  - allow more than 2 ordinates in each dimension
- *E.g.* mesh (2-cube), 3D-mesh (3-cube)
- Matches with physical world structure
- Bounds degree at node
- Has Locality
- **Even more bottleneck potentials**
  - make channels wider (CS184a:Day 17)

Caltech CS184 Spring2005 -- DeHon

19

## Torus

- Wrap around n-cube ends
  - 2-cube → cylinder
  - 3-cube → donut
- Cuts worst-case distances in half
- Can be laid-out reasonable efficiently
  - maybe 2x cost in channel width?



Caltech CS184 Spring2005 -- DeHon

20

## Fat-Tree

- Saw that communications typically has locality (CS184a)
- Modeled recursive bisection/Rent's Rule
- Leiserson showed Fat-Tree was (area, volume) universal
  - w/in  $\log(N)$  the area of **any** other structure
  - exploit physical space limitations wiring in {2,3}-dimensions

Caltech CS184 Spring2005 -- DeHon

21

## MoT/Express Cube (Mesh with Bypass)

- Large machine in 2 or 3 D mesh
  - routes must go through square/cube root switches
  - vs.  $\log(N)$  in fat-tree, hypercube, MIN
- Saw practically can go further than one hop on wire...
- Add long-wire bypass paths

Caltech CS184 Spring2005 -- DeHon

22

## Routing Styles

Caltech CS184 Spring2005 -- DeHon

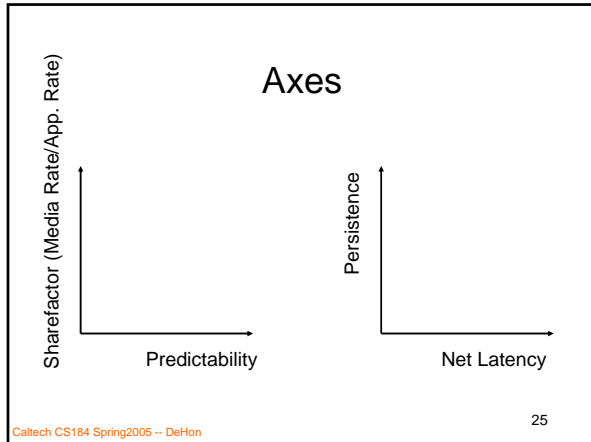
23

## Issues/Axes

- Throughput of Communication relative to data rate of media
  - Single point-to-point link consume media BW?
  - Can share links between multiple comm streams?
  - What is the sharing factor?
- Binding time/Predictability of Interconnect
  - Pre-fab
  - Before communication then use for long time
  - Cycle-by-cycle
- Network latency vs. persistence of communication
  - Comm link persistence

Caltech CS184 Spring2005 -- DeHon

24



### Hardwired

- Direct, fixed wire between two points
- *E.g.* Conventional gate-array, std. cell
- Efficient when:
  - know communication *a priori*
    - fixed or limited function systems
    - high load of fixed communication
      - often control in general-purpose systems
  - links carry high throughput traffic continually between fixed points

Caltech CS184 Spring2005 -- DeHon

26

### Configurable

- Offline, lock down persistent route.
- *E.g.* FPGAs
- Efficient when:
  - link carries high throughput traffic
    - (loaded usefully near capacity)
  - traffic patterns change
    - on timescale  $\gg$  data transmission

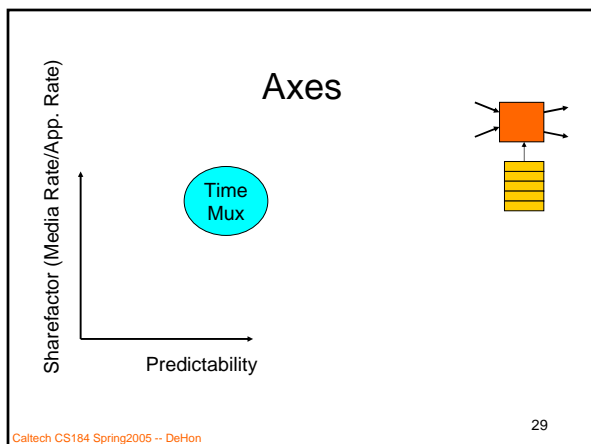
Caltech CS184 Spring2005 -- DeHon

### Time-Switched

- Statically scheduled, wire/switch sharing
- *E.g.* TDMA, NuMesh, TSFPGA
- Efficient when:
  - thput per channel  $<$  thput capacity of wires and switches
  - traffic patterns change
    - on timescale  $\gg$  data transmission

Caltech CS184 Spring2005 -- DeHon

28

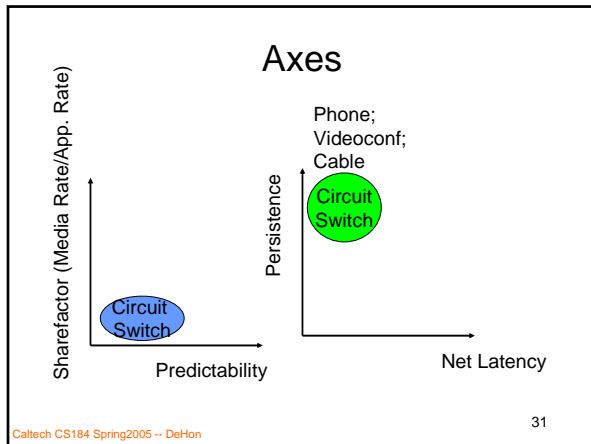


### Self-Route, Circuit-Switched

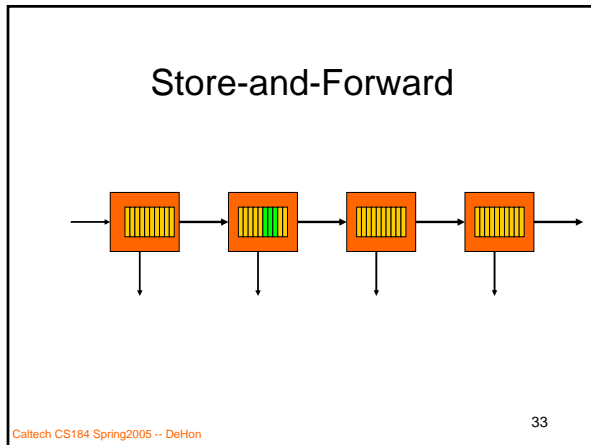
- Dynamic arbitration/allocation, lock down routes
- *E.g.* METRO/RN1
- Efficient when:
  - instantaneous communication bandwidth is high (consume channel)
  - lifetime of comm.  $>$  delay through network
  - communication pattern unpredictable
  - rapid connection setup important

Caltech CS184 Spring2005 -- DeHon

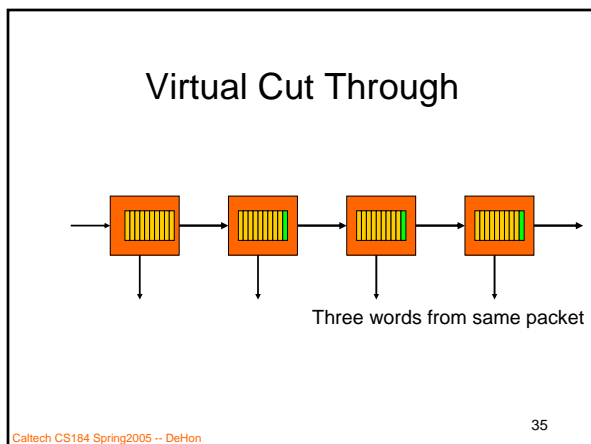
30



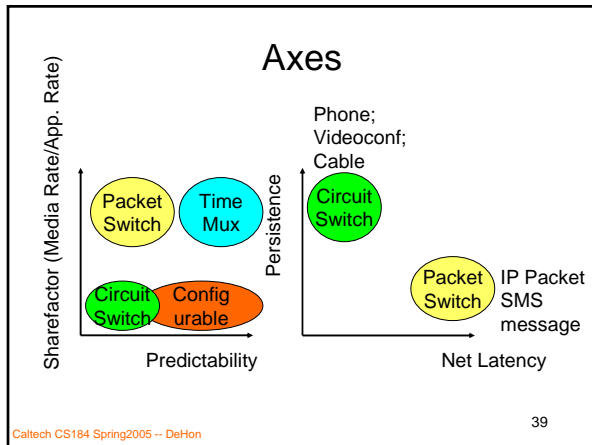
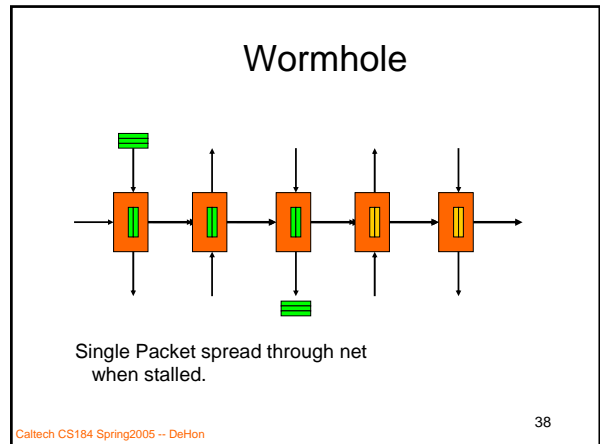
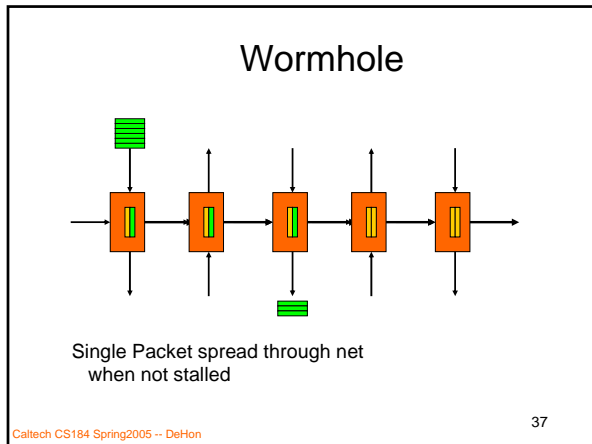
- ### Self-Route, Store-and-Forward, Packet Switched
- Dynamic arbitration, packetized data
  - Get entire packet before sending to next node
  - E.g. nCube, early Internet routers
  - Efficient when:
    - lifetime of comm < delay through net
    - communication pattern unpredictable
    - can provide buffer/consumption guarantees
    - packets small
- 32
- Caltech CS184 Spring2005 -- DeHon



- ### Self-Route, Virtual Cut Through
- Dynamic arbitration, packetized data
  - Start forwarding to next node as soon as have header
  - Don't pay full latency of storing packet
  - Keep space to buffer entire packet if necessary
  - Efficient when:
    - lifetime of comm < delay through net
    - communication pattern unpredictable
    - can provide buffer/consumption guarantees
    - packets small
- 34
- Caltech CS184 Spring2005 -- DeHon



- ### Self-Route, Wormhole Packet-Switched
- Dynamic arbitration, packetized data
  - E.g. Caltech MRC, Modern Internet Routers
  - Efficient when:
    - lifetime of comm < delay through net
    - communication pattern unpredictable
    - can provide buffer/consumption guarantees
    - message > buffer length
      - allow variable (? Long) sized messages
- 36
- Caltech CS184 Spring2005 -- DeHon



## Online Routing

Caltech CS184 Spring2005 -- DeHon 40

- ### Costs: Area
- Area
    - switch ( $1-1.5K\lambda^2$ / switch)
      - larger with pipeline ( $4K\lambda^2$ ) and rebuffer
    - state (SRAM bit =  $1.2K\lambda^2$ / bit)
      - multiple in time-switched cases
    - arbitration/decision making
      - usually dominates above
    - buffering (SRAM cell per buffer)
      - can dominate
- Time Mux  
Dynamic  
Caltech CS184 Spring2005 -- DeHon 41

### Area

$$A_{switching} = W \times Ins \times Outs \times A_{sw}$$

$$A_{inst} = Outs \times \log_2(Ins) \times Cycles \times A_{bit}$$

$$A_{arb} = Ins \times Outs \times A_{arbitration}$$

$$A_{queue} = W \times (Ins + Outs) \times Depth \times A_{bit}$$

(queue rough approx; you will refine) 42  
Caltech CS184 Spring2005 -- DeHon

## Costs: Latency

- Time single path
  - make decisions
  - round-trip flow-control
- Time contention/traffic
  - blocking in buffers
  - quality of decision
    - pick wrong path
    - have stale data

Caltech CS184 Spring2005 -- DeHon

43

## Intermediate Approach

- For large # of predictable patterns
  - switching memory may dominate allocation area
  - area of routed case < time-switched
  - [e.g. Large Cycles]
- Get offline, global planning advantage
  - by **source routing**
    - source specifies offline determined route path
    - offline plan avoids contention

Caltech CS184 Spring2005 -- DeHon

44

## Offline vs. Online

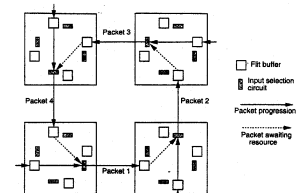
- If know patterns in advance
  - offline cheaper
    - no arbitration (area, time)
    - no buffering
    - use more global data
      - better results
- As becomes less predictable
  - benefit to online routing

Caltech CS184 Spring2005 -- DeHon

45

## Deadlock

- Possible to introduce deadlock
- Consider wormhole routed mesh



[example from Li and McKinley, IEEE Computer v26n2, 1993]

Caltech CS184 Spring2005 -- DeHon

46

## Dimension Order Routing

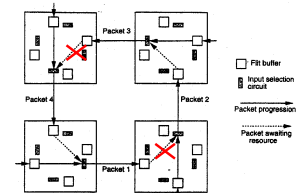
- Simple (early Caltech) solution
  - order dimensions
  - force complete routing in lower dimensions before route in next higher dimension

Caltech CS184 Spring2005 -- DeHon

47

## Dimension Ordered Routing

- Route Y, then Route X



[example from Li and McKinley, IEEE Computer v26n2, 1993]

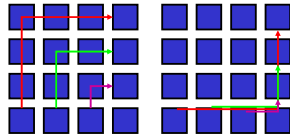
Caltech CS184 Spring2005 -- DeHon

48



## Dimension Order Routing

- Avoids cycles in channel graph
- Limits routing freedom
- Can cause artificial congestion
  - consider
    - (0,0) to (3,3)
    - (1,0) to (3,2)
    - (2,0) to (3,1)



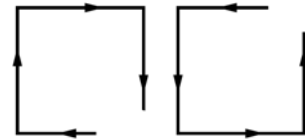
• [There is a rich literature on how to do better]

Caltech CS184 Spring2005 -- DeHon

49

## Turn Model

- Problem is cycles
- Selectively disallow turns to break cycles
- 2D Mesh



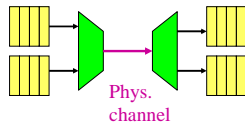
West-First Routing

Caltech CS184 Spring2005 -- DeHon

50

## Virtual Channel

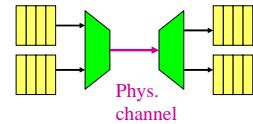
- **Variation:** each physical channel represents multiple logical channels
  - each logical channel has own buffers
  - blocking in one VC allows other VCs to use the physical link



Caltech CS184 Spring2005 -- DeHon

51

## Virtual Channels



- Benefits
  - can be used to remove cycles
    - e.g. separate increasing and decreasing channels
    - route increasing first, then decreasing
    - more freedom than dimension ordered
  - prioritize traffic
    - e.g. prevent control/OS traffic from being blocked by user traffic
  - better utilization of physical routing channels

Caltech CS184 Spring2005 -- DeHon

52

## Lost Freedom?

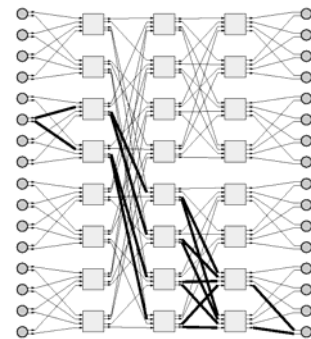
- Online routes often make (must make) decisions based on local information
- Can make wrong decision
  - i.e. two paths look equally good at one point in net
    - but one leads to congestion/blocking further ahead

Caltech CS184 Spring2005 -- DeHon

53

## Multibutterfly Network

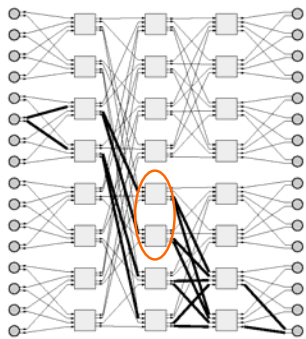
- Dilated routers
  - have multiple outputs in each logical direction
  - Means multiple paths between any src, sink pair
- Use to avoid congestion
  - also faults



Caltech CS184 Spring2005 -- DeHon

## Multibutterfly Network

- Can get into local blocking when there is a path
- Costs of not having global information



Caltech CS184 Spring2005 -- DeHon

## Transit/Metro

- Self-routing circuit switched network
- When have choice
  - select randomly
    - avoid bad structural cases
- When blocked
  - drop connection
  - allow to route again from source
  - stochastic search explores all paths
    - finds any available

Caltech CS184 Spring2005 -- DeHon

56

## Relation to Project

Caltech CS184 Spring2005 -- DeHon

57

## Intuitive Tradeoff

- Benefit of Time-Multiplexing?
  - Minimum end-to-end latency
  - No added decision latency at runtime
  - Offline route → high quality route
    - → use wires efficiently
- Cost of Time-Multiplexing?
  - Route task must be static
    - Cannot exploit low activity
  - Need memory bit per switch per time step
    - Lots of memory if need large number of time steps...

Caltech CS184 Spring2005 -- DeHon

58

## Intuitive Tradeoff

- Benefit of Packet Switching?
  - No area proportional to time steps
  - Route only active connections
  - Avoids slow, off-line routing
- Cost of Packet Switching?
  - Online decision making
    - Maybe won't use wires as well
    - Potentially slower routing?
      - Slower clock, more clocks across net
  - Data will be blocked in network
    - Adds latency
    - Requires packet queues

Caltech CS184 Spring2005 -- DeHon

59

## Packet Switch Motivations

- SMVM:
  - Long offline routing time limits applicability
  - Route memory exceeds compute memory for large matrices
- ConceptNet:
  - Evidence of low activity for keyword retrieval ... could be important to exploit

Caltech CS184 Spring2005 -- DeHon

60

## Example

- ConceptNet retrieval
  - Visits 84K nodes across all time steps
  - 150K nodes
  - 8 steps  $\rightarrow$  1.2M node visits
  - Activity less than 7%

Caltech CS184 Spring2005 -- DeHon

61

## Dishoom/ConceptNet Estimates

$N_z$	$N_{FPGA}$	$T_{comp}$	$T_{load}$	$T_{lat}$	$T_{step}$
1	64	500	1500	48	2048
2	128	250	750	52	1052
4	256	125	375	60	560
8	512	63	188	76	327

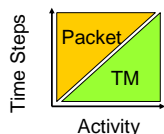
Pushing all nodes, all edges;  
Bandwidth ( $T_{load}$ ) dominates.

Caltech CS184 Spring2005 -- DeHon

62

## Question

- For what activity factor does Packet Switching beat Time Multiplexed Routing?
  - To what extent is this also a function of total time steps?



Caltech CS184 Spring2005 -- DeHon

63

## Admin

- Reading
- VHDL intro on Wednesday
- Fast Virtex Queue Implementations on Friday

Caltech CS184 Spring2005 -- DeHon

64

## Big Ideas

- Must work with constraints of physical world
  - only have 3 dimensions (2 on current VLSI) in which to build interconnect
  - Interconnect can be dominate area, time
  - gives rise to universal networks
    - e.g. fat-tree

Caltech CS184 Spring2005 -- DeHon

65

## Big Ideas

- Structure
  - exploit physical locality where possible
  - the more predictable behavior
    - cheaper the solution
  - exploit earlier binding time
    - cheaper configured solutions
    - allow higher quality offline solutions
- Interconnect style
  - Driven by technology and application traffic patterns

Caltech CS184 Spring2005 -- DeHon

66