

## Interconnection Network Topology Design Trade-offs

CS 258, Spring 99  
 David E. Culler  
 Computer Science Division  
 U.C. Berkeley

## Real Machines

Machine	Topology	Cycle Time (ns)	Channel Width (bits)	Routing Delay (cycles)	Rat. (data/ops)
eCUBE2	Hypercube	25	1	40	32
TMC CM-5	Fat-tree	25	4	10	4
IBM SP-2	Bus	25	8	5	16
Intel Paragon	2D Mesh	11.5	16	2	16
Melba C3-2	Fat-Tree	20	8	7	8
Cray T3E	3D Torus	6.67	16	2	16
DASH	Torus	30	16	2	16
VMachine	3D Mesh	31	8	2	8
Mosaic	Butterfly	20	16	2	16
SGI Origin	Hypercube	2.5	20	18	180
Myrinet	Arbitrary	6.25	16	50	16

- Wide links, smaller routing delay
- Tremendous variation

## Interconnection Topologies

- Class networks scaling with N
- Logical Properties:
  - distance, degree
- Physical properties
  - length, width
- Fully connected network
  - diameter = 1
  - degree = N
  - cost?
    - bus => O(N), but BW is O(1) - actually worse
    - crossbar => O(N<sup>2</sup>) for BW O(N)
- VLSI technology determines switch degree

## Linear Arrays and Rings

- Linear Array
  - Diameter?
  - Average Distance?
  - Bisection bandwidth?
  - Route A -> B given by relative address R = B - A
- Torus?
  - Examples: FDDI, SCI, FiberChannel Arbitrated Loop, KSR1

## Multidimensional Meshes and Tori

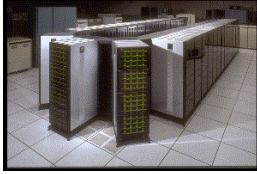
- d-dimensional array**
  - n = k<sub>d-1</sub> X ... X k<sub>0</sub> nodes
  - described by d-vector of coordinates (i<sub>d-1</sub>, ..., i<sub>0</sub>)
- d-dimensional k-ary mesh: N = k<sup>d</sup>**
  - k = dN
  - described by d-vector of radix k coordinate
- d-dimensional k-ary torus (or k-ary d-cube)?**

## Properties

- Routing**
  - relative distance: R = (b<sub>d-1</sub> - a<sub>d-1</sub>, ..., b<sub>0</sub> - a<sub>0</sub>)
  - traverse r<sub>i</sub> = b<sub>i</sub> - a<sub>i</sub> hops in each dimension
  - dimension-order routing
- Average Distance**
  - d x 2k/3 for mesh
  - dk/2 for cube
- Degree?**
- Bisection bandwidth?**
  - k<sup>d-1</sup> bidirectional links
- Physical layout?**
  - 2D in O(N) space
  - higher dimension?

Wire Length?  
Partitioning?  
Short wires

## Real World 2D mesh



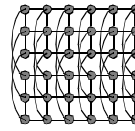
- 1824 node Paragon: 16 x 114 array

3/19/99

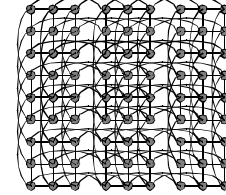
CS258 S99

7

## Embeddings in two dimensions



6 x 3 x 2



- Embed multiple logical dimension in one physical dimension using long wires

3/19/99

CS258 S99

8

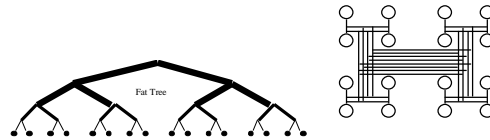
## Trees



- Diameter and ave distance logarithmic
  - k-ary tree, height  $d = \log_k N$
  - address specified d-vector of radix k coordinates describing path down from root
- Fixed degree
- Route up to common ancestor and down
  - $R = B \text{ xor } A$
  - let  $i$  be position of most significant 1 in  $R$ , route up  $i+1$  levels
  - down in direction given by low  $i+1$  bits of  $B$
- H-tree space is  $O(N)$  with  $O(\sqrt{N})$  long wires
- Bisection BW? CS258 S99

9

## Fat-Trees



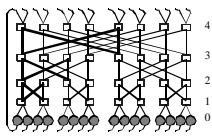
- Fatter links (really more of them) as you go up, so bisection BW scales with  $N$

3/19/99

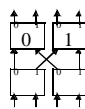
CS258 S99

10

## Butterflies



16 node butterfly



building block

- Tree with lots of roots!
- $N \log N$  (actually  $N/2 \times \log N$ )
- Exactly one route from any source to any dest
- $R = A \text{ xor } B$ , at level  $i$  use 'straight' edge if  $r_i=0$ , otherwise cross edge
- Bisection  $N/2$  vs  $n^{(d-1)/d}$  CS258 S99

3/19/99

CS258 S99

11

## k-ary d-cubes vs d-ary k-flies

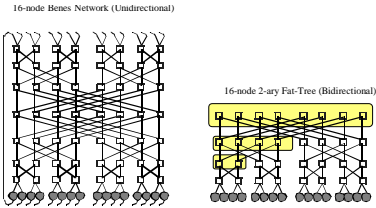
- degree  $d$
- $N$  switches vs  $N \log N$  switches
- diminishing BW per node vs constant
- requires locality vs little benefit to locality
- Can you route all permutations?

3/19/99

CS258 S99

12

## Benes network and Fat Tree



- Back-to-back butterfly can route all permutations – off line
- What if you just pick a random mid point?

3/19/99

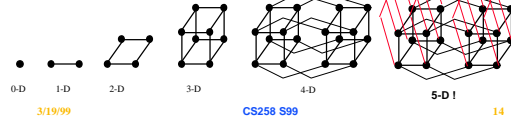
CS258 S99

13

## Hypercubes

- Also called binary n-cubes. # of nodes =  $N = 2^n$ .
- $O(\log N)$  Hops
- Good bisection BW
- Complexity
  - Out degree is  $n = \log N$

correct dimensions in order  
– with random comm. 2 ports per processor

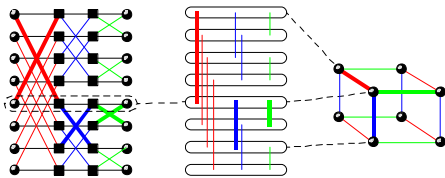


3/19/99

CS258 S99

14

## Relationship BttrFlies to Hypercubes



- Wiring is isomorphic
- Except that Butterfly always takes  $\log n$  steps

3/19/99

CS258 S99

15

## Topology Summary

Topology	Degree	Diameter	Ave Dist	Bisection	D (D ave) @ P=1024
1D Array	2	$N-1$	$N/3$	1	huge
1D Ring	2	$N/2$	$N/4$	2	
2D Mesh	4	$2(N^{1/2} - 1)$	$2/3 N^{1/2}$	$N^{1/2}$	63 (21)
2D Torus	4	$N^{1/2}$	$1/2 N^{1/2}$	$2N^{1/2}$	32 (16)
k-ary n-cube	$2n$	$nk/2$	$nk/4$	$nk/4$	15 (7.5) @n=3
Hypercube	$n = \log N$	$n$	$n$	$n/2$	10 (5)

- All have some “bad permutations”
  - many popular permutations are very bad for meshes (transpose)
  - randomness in wiring or routing makes it hard to find a bad one!

3/19/99

CS258 S99

16

## How Many Dimensions?

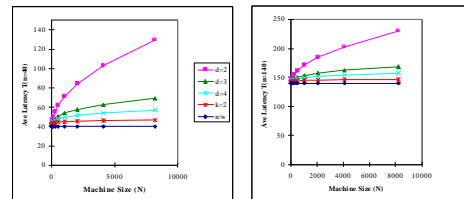
- $n = 2$  or  $n = 3$ 
  - Short wires, easy to build
  - Many hops, low bisection bandwidth
  - Requires traffic locality
- $n \geq 4$ 
  - Harder to build, more wires, longer average length
  - Fewer hops, better bisection bandwidth
  - Can handle non-local traffic
- k-ary d-cubes provide a consistent framework for comparison
  - $N = k^d$
  - scale dimension (d) or nodes per dimension (k)
  - assume cut-through

3/19/99

CS258 S99

17

## Traditional Scaling: Latency(P)



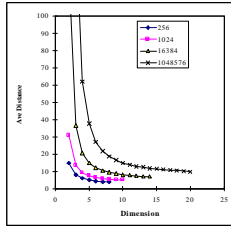
- Assumes equal channel width
  - independent of node count or dimension
  - dominated by average distance

3/19/99

CS258 S99

18

## Average Distance



$$\text{ave dist} = d(k-1)/2$$

- but, equal channel width is not equal cost!
- Higher dimension => more channels

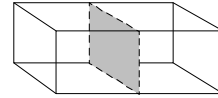
3/19/99

CS258 S99

19

## In the 3D world

- For  $n$  nodes, bisection area is  $O(n^{2/3})$



- For large  $n$ , bisection bandwidth is limited to  $O(n^{2/3})$ 
  - Bill Dally, IEEE TPDS, [Dal90a]
  - For fixed bisection bandwidth, low-dimensional  $k$ -ary  $n$ -cubes are better (otherwise higher is better)
  - i.e., a few short fat wires are better than many long thin wires
  - What about many long fat wires?

3/19/99

CS258 S99

20

## Equal cost in $k$ -ary $n$ -cubes

- Equal number of nodes?
- Equal number of pins/wires?
- Equal bisection bandwidth?
- Equal area? Equal wire length?

What do we know?

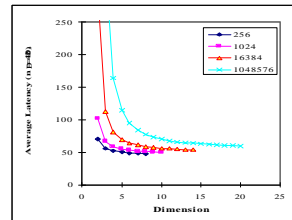
- switch degree:  $d$  diameter =  $d(k-1)$
- total links =  $Nd$
- pins per node =  $2wd$
- bisection =  $k^{d-1} = N/k$  links in each direction
- $2Nw/k$  wires cross the middle

3/19/99

CS258 S99

21

## Latency( $d$ ) for $P$ with Equal Width



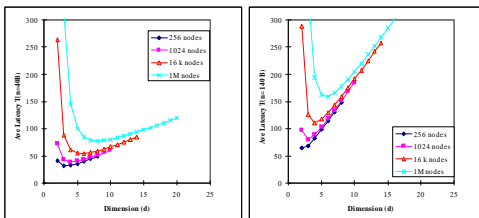
- total links( $N$ ) =  $Nd$

3/19/99

CS258 S99

22

## Latency with Equal Pin Count



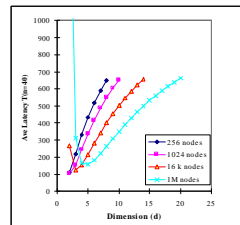
- Baseline  $d=2$ , has  $w = 32$  (128 wires per node)
- fix  $2dw$  pins =>  $w(d) = 64/d$
- distance up with  $d$ , but channel time down

3/19/99

CS258 S99

23

## Latency with Equal Bisection Width



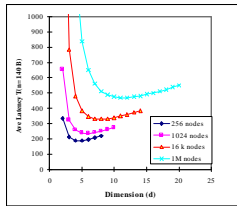
- $N$ -node hypercube has  $N$  bisection links
- 2d torus has  $2N^{1/2}$
- Fixed bisection =>  $w(d) = N^{1/d} / 2 = k/2$
- 1 M nodes,  $d=2$  has  $w=512!$

3/19/99

CS258 S99

24

## Larger Routing Delay (w/ equal pin)



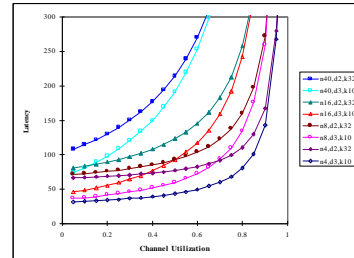
- Dally's conclusions strongly influenced by assumption of small routing delay

3/19/99

CS258 S99

25

## Latency under Contention



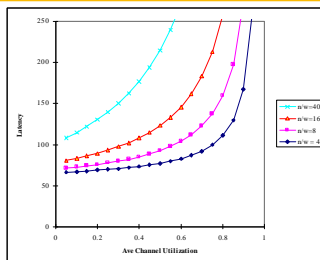
- Optimal packet size? Channel utilization?

3/19/99

CS258 S99

26

## Saturation



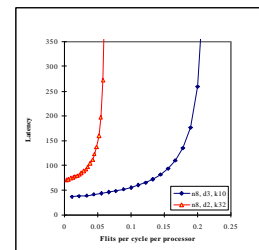
- Fatter links shorten queuing delays

3/19/99

CS258 S99

27

## Phits per cycle



- higher degree network has larger available bandwidth

3/19/99— cost?

CS258 S99

28

## Discussion

- Rich set of topological alternatives with deep relationships
- Design point depends heavily on cost model
  - nodes, pins, area, ...
  - Wire length or wire delay metrics favor small dimension
  - Long (pipelined) links increase optimal dimension
- Need a consistent framework and analysis to separate opinion from design
- Optimal point changes with technology

3/19/99

CS258 S99

29