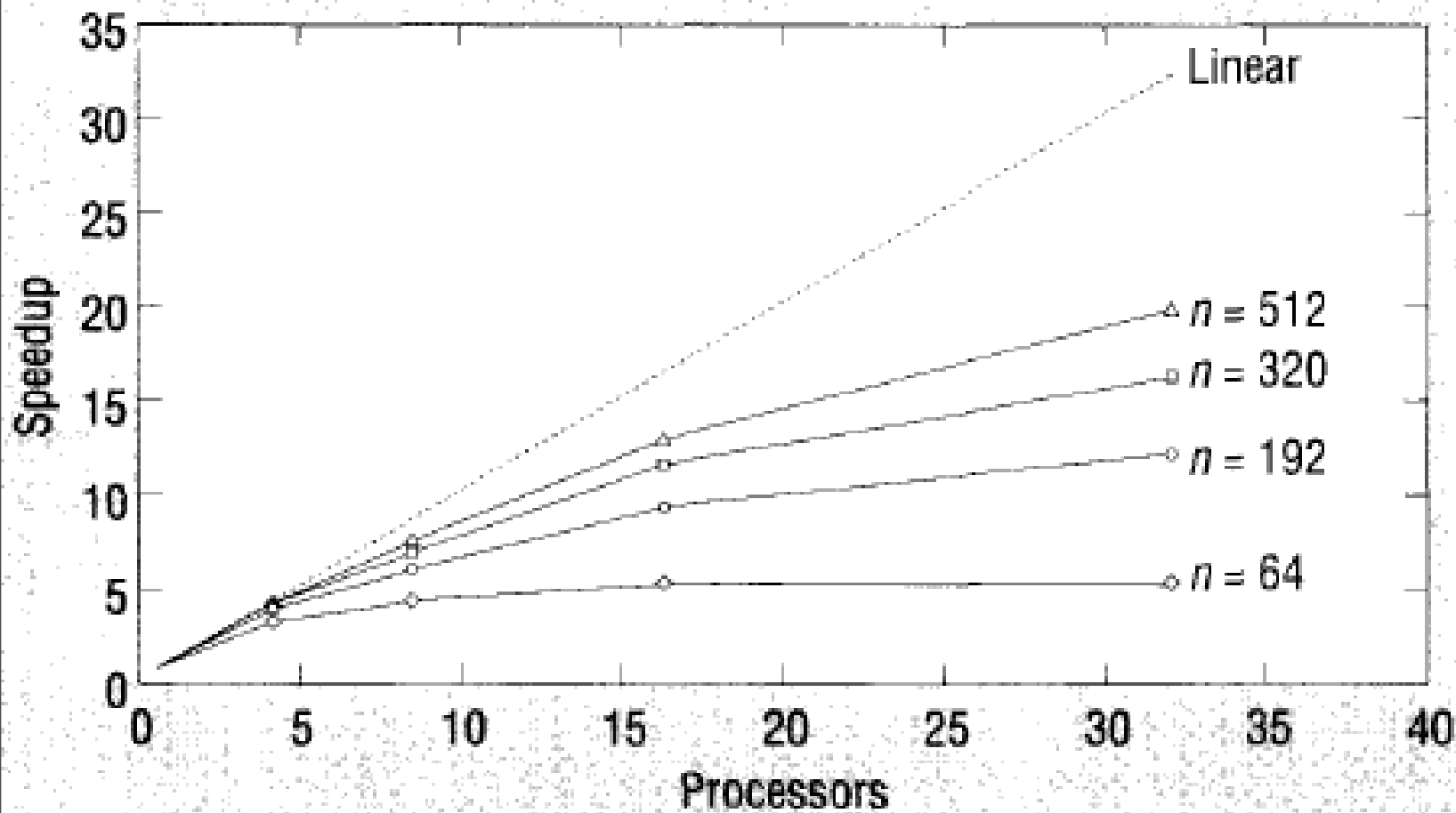


Isoefficiency analysis

Measuring the parallel scalability of algorithms

- One of many parallel performance metrics
- Allows us to determine scalability with respect to machine parameters
 - number of processors and their speed
 - communication patterns, bandwidth and startup
- Give us a way of computing
 - the relative scalability of two algorithms
 - how much work needs to be increased when the number of processors is increased to maintain the same efficiency

Amdahl's law reviewed



As number of processors increase, serial overheads reduce efficiency

As problem size increases, efficiency returns

Efficiency of adding n numbers on an ancient machine

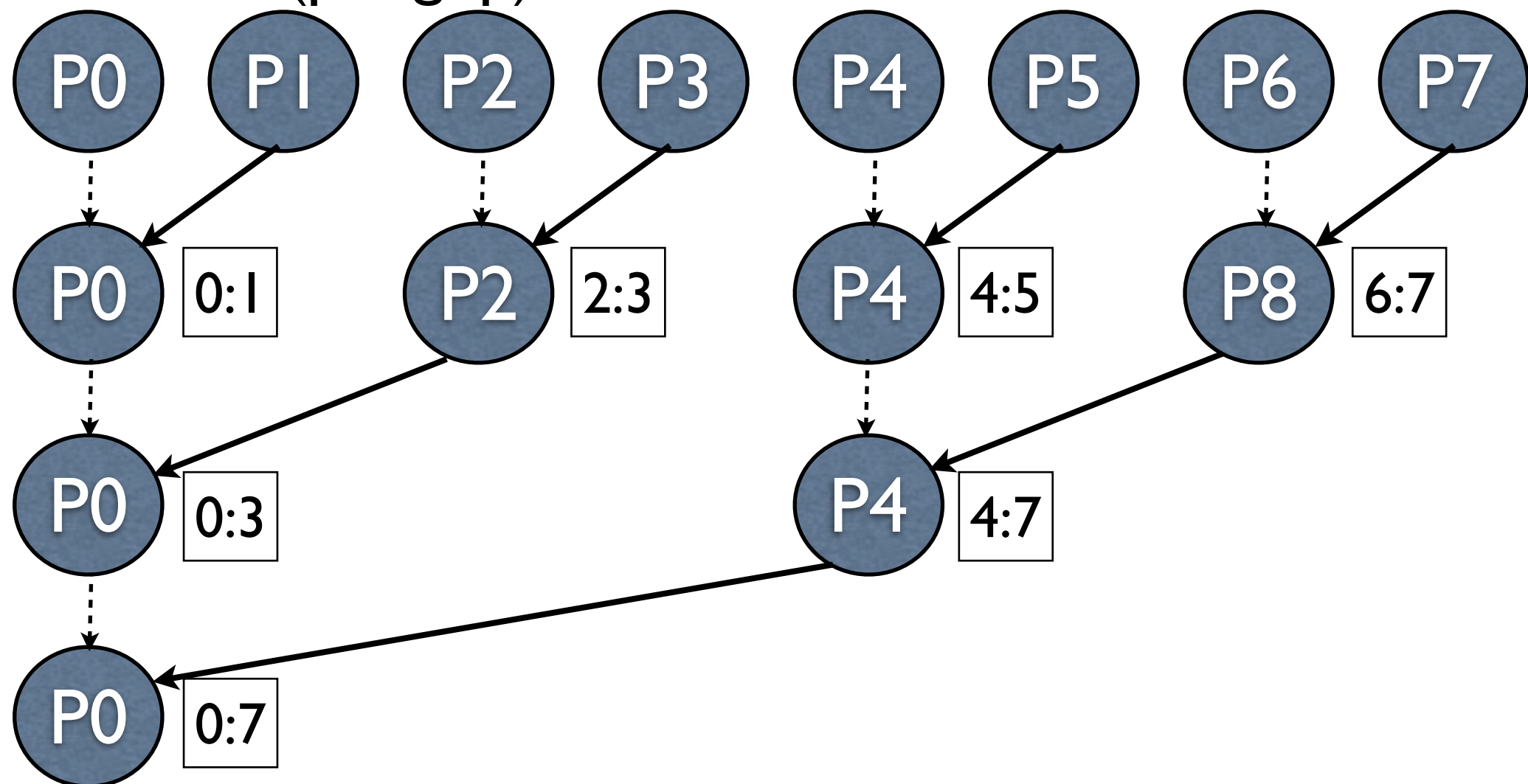
$P=4$ gives ϵ of .80 with 64 numbers

$P=8$ gives ϵ of .80 with 192 numbers

$P=16$ gives ϵ of .80 with 512 numbers

Motivating example

- Consider a program that does $O(n)$ work
- Also assume the overhead is $O(\log_2 p)$, i.e. it does a reduction
- The *total overhead*, i.e. the amount of time processors are sitting idle or doing work associated with parallelism instead of the basic problem, is $O(p \log_2 p)$



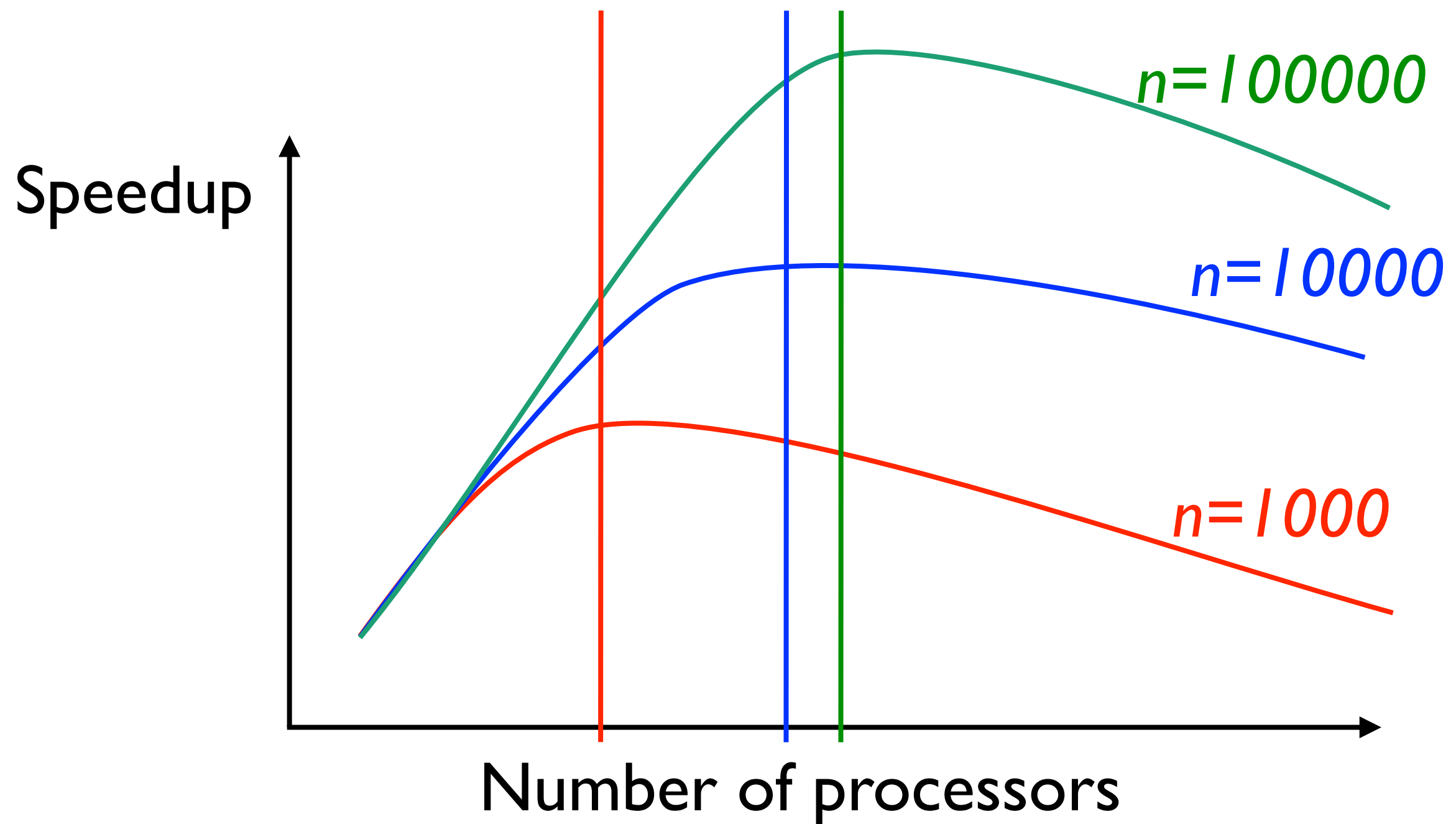
Data to maintain efficiency

P	$P \log_2 P$	Data needed
2	2	1 GB
4	8	4 GB
8	24	12 GB
16	64	32

- What if the machine only has 12 GB of memory?
- Efficiency will drop off after 8 processors
- Cannot increase data and problem size fast enough to overcome parallel overheads
- If work done per datum increased, could scale further

Isoefficiency analysis allows us to analyze the rate at which the data size must grow to mask parallel overheads to determine if a computation is scalable.

Amdahl Effect both increases speedup and move the knee of the curve to the right



Total overhead T_o is the time spent

- Any time not spent by the serial form of the program
- In communication
- Idle time because of waiting for some processor that is executing serial code
- Idle time waiting for data from another processor
- ...

Efficiency revisited

- Total time spent on all processors is the original sequential execution [*best sequential implementation*] time plus the parallel overhead

$$PT_p = T_1 + T_O$$

- The time it takes the program to run in parallel is the total time spent on all processors divided by the number of processors. This is true because T_O includes the time processors are waiting for something to happen in a parallel execution.

$$T_p = (T_1 + T_O)/P \quad (1), \text{ which can be written } T_1 = P T_p - T_O \quad (1a)$$

- Speedup is as before (T_1/T_P), or by substituting (1, 1a) above, we get:

$$\begin{aligned} S &= (P T_p - T_O) / ((T_1 + T_O)/P) = (P^2 T_p - P T_O) / (T_1 + T_O) \text{ using 1 we get} \\ &= P(T_1 + T_O) - P T_O / (T_1 + T_O) = P (T_1 + T_O - T_O) / (T_1 + T_O) \\ &= P T_1 / (T_1 + T_O) \end{aligned}$$

Efficiency revisited

- With speedup being

$$S = T_1/T_P = (P T_1)/(T_1 + T_o)$$

- Efficiency can be computed using the previous definition of the ratio of S to P as:

$$\begin{aligned} E = S/P &= ((P T_1)/(T_1 + T_o))/P = T_1/(T_1 + T_o) \\ &= 1/(1 + T_o/T_1) \end{aligned}$$

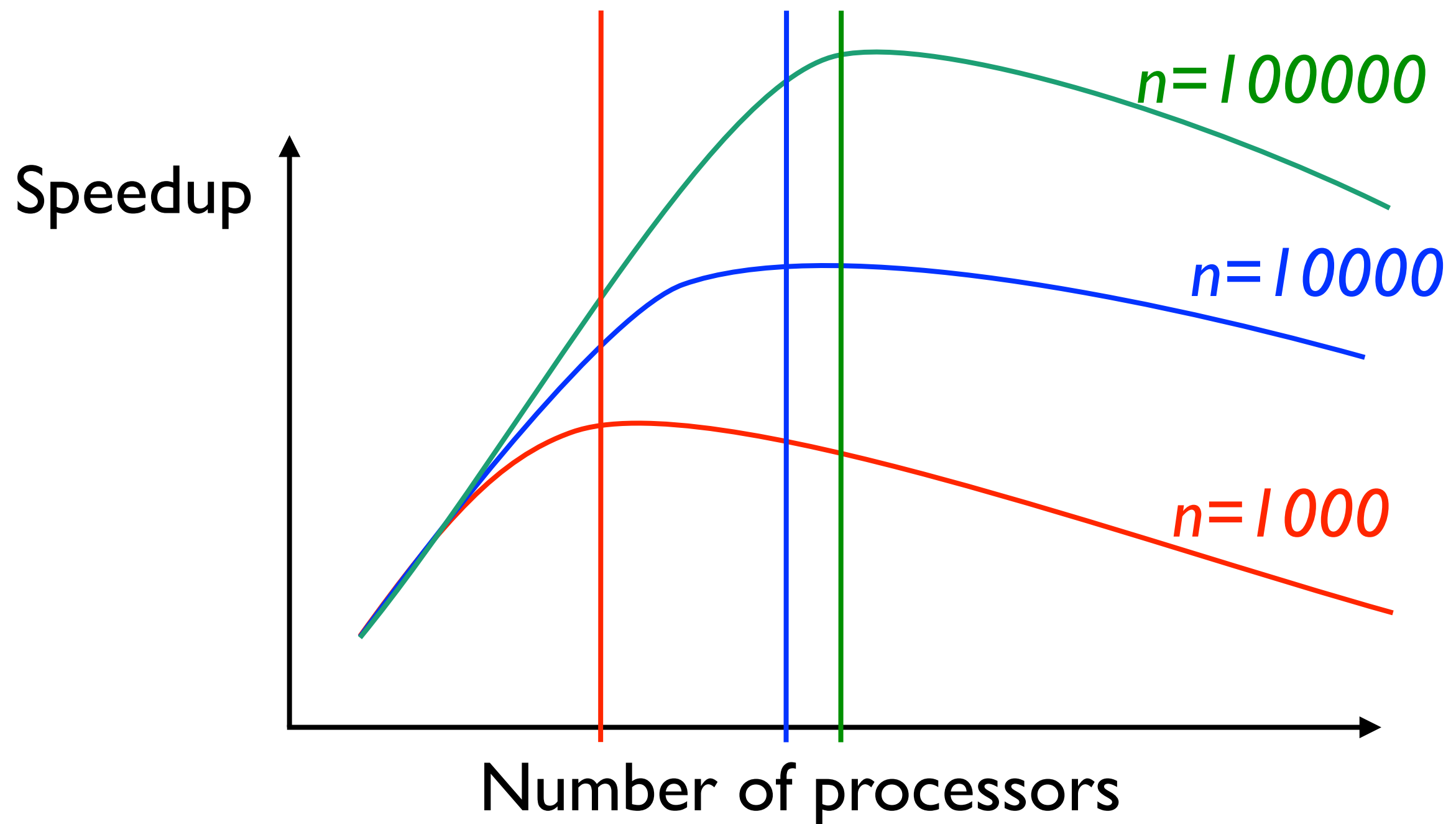
Efficiency as a function of work, data size and overhead

- Let
 - T_1 be the single processor time
 - W be the amount of work of units of work)
 - t_c be the time to perform each unit of work
- Then $T_1 = W \cdot t_c$
- T_O is the total overhead, i.e. time spent doing parallel stuff but not the original work
- Then *efficiency* can be rewritten as $E = \frac{1}{1 + \frac{T_o}{W t_c}}$

Some insights into Amdahl's law and the Amdahl effect can be gleaned from this

- Efficiency is $E = \frac{1}{1 + \frac{T_o}{Wt_c}}$
- For the same problem size on more processors, W is constant and T_o is growing. **Thus efficiency decreases.**
- Let $\theta(W)$ be all functions that grow at the same or **faster rate than** W , i.e. $\theta(W)$ is an upper bound
- As P increases, T_o will grow faster, the same, or slower than $\theta(W)$
 - If faster, system has limited scalability
 - If slower or the same, system is very scalable, can grow work the same or slower than processor growth

Amdahl Effect both increases speedup and move the knee of the curve to the right



The relationship of work and constant efficiency

$$E = \frac{1}{1 + \frac{T_o}{Wt_c}}$$

$$\frac{T_o}{W} = t_c \left(\frac{1 - E}{E} \right)$$

$$W = \frac{1}{t_c} \left(\frac{E}{1 - E} \right) T_o$$

If $K = E/(t_c(1 - E))$ is a constant that depends on the efficiency, then we can reduce the last equation to

$$W = KT_o$$


This relationship holds when the efficiency is constant

Will use algebraic manipulations to (eventually) represent W as a function of P . This indicates how W must grow as the number of processors grows to maintain the same efficiency.

What if T_0 negative?

- *Superlinear* speedups can lead to negative values for T_0
- Appears to cause work to need to decrease
- Causes of *superlinear* speedup
 - increased memory size in NUMA and hierarchical (i.e. caching) memory systems
 - Search based problems
- We assume $T_0 \geq 0$

Simple case superlinear speedup

-- linear scan search for 9

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----

To find an element takes $O(\text{pos})$ steps

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----

To find an element takes $O(\text{pos} - \text{offset})$ steps

Doubling processors leads to a speedup of ~ 9

Cache effects

- Moving data into cache is a hidden work item
- As P grows larger, total cache available also grows larger
- If P grows faster than data size, eventually all data fits in cache
- Thus cache misses due to capacity vanish, and the work associated with those misses vanishes, and the parallel program is doing less work, enabling superlinear speedups if everything else is highly efficient

There are no magical causes of superlinear speedup

- In the early days of parallel computing it was thought by some that there might be something else going on with parallel executions
- All cases of superlinear speedup observed to date can be explained by a reduction in overall work in solving the problem

How to define the problem size or amount of work

- Given an $n \times n$ matrix multiply, what is the problem size?
 - Commonly called n
- How about adding two $n \times n$ matrices?
 - Could call it the same n
- How about adding to vectors of length n ?
 - Could also call the problem size n
- Yet one involves n^3 work, one n^2 work, and one n work

Same name, different amounts of work - this causes problems

- The goal of isoefficiency is to see how work should scale to maintain efficiency
- Let $W=n$ for matrix multiply, matrix add and vector addition
- Let all three (for the sake of this example, *even though not true*) have a similar T_0 that grows linearly with P
- Doubling n would lead to 8 times more operations for matrix multiply, 4 times for matrix add, and 2 times for vector add
- Intuitively the vector add seems to be right, since number of operations and work (W) seem to be the same thing, not data size.
- We will normalize W in terms of *operations in the best sequential algorithm*, not some other metric of problem size

Adding n numbers

- $n-1$ operations in sequential algorithm -- asymptotically is n and we will use n , and $T_1 = n \cdot t_c$
- Let each add take one unit of time, and each communication take one unit of time
- On P processors, n/P operations + $\log_2 P$ communication steps + one add operation at each communication step
- $T_P = n/P + 2 \log_2 P$
- $T_O = P (2 \log_2 P)$ since each processor is either doing this or waiting for this to finish on other processors
- $S = T_1 / T_P = n / (n/P + 2 \log_2 P)$
- $E = S / P = n / (n + 2 P \log_2 P)$

Isoefficiency analysis of adding n numbers

- From slide 12, $W = K T_O$ if same efficiency is to be maintained
- $T_O = P (2 \log_2 P)$ from the previous slide, then

$$W = 2 K P \log_2 P$$

and ignoring constants give an isoefficiency function of

$$\theta(P \log_2 P)$$

- If the number of processors is increased to P' , then the work must be increased not by P'/P , but by
 $(P' \log_2 P') / (P \log_2 P)$
- Thus going from 4 to 16 processors requires having
 $(16 \log_2 16) / (4 \log_2 4)$ or 8 times as much work

More complicated T_O

- Consider $T_O = P^{3/2} + P^{3/4}W^{3/4}$, and $W = K T_O$, then
 $W = P^{3/2} + P^{3/4}W^{3/4}$ (again, ignoring constant K)
- Difficult to solve for W in terms of P
- Note that we need ratio of W and T_O to remain fixed for E (efficiency) to remain fixed
- Problem will scale well if no term of T_O grows faster than W
- Thus we can examine terms independently

$$W = P^{3/2} + P^{3/4} W^{3/4}$$

- Solve for first term, i.e. $W = KP^{3/2} = \theta(P^{3/2})$
- Solve for second term, i.e.

$$W = K P^{3/4} W^{3/4}$$

$$W^{1/4} = K P^{3/4}$$

$$W = K^4 P^3 = \theta(P^3)$$

- If problem size grows at least as fast as $\theta(P^{3/2})$ and $\theta(P^3)$ then efficiency will not decrease as P increases.
- Thus the isoefficiency function for the system is $\theta(P^3)$

Cost optimality

- Parallel system is *cost-optimal* if product of $PT_P \propto W$
i.e, is not growing faster than W
- Because $PT_p = T_1 + T_O$, then $T_1 + T_O \propto W$
- Since $T_1 = Wt_c$, we have $Wt_c + T_O \propto W$ and therefore $W \propto T_O$.
- Suggests a parallel system is cost optimal if its overhead function and problem size are of the same order of magnitude, i.e. have same order of complexity.
- **Conforming to the isoefficiency relationship keeps a system cost-optimal as it is scaled up**

How small can isoefficiency function be?

- Let a problem contain W basic operations
- Let problem size grow slower than $\theta(P)$
- As P grows, eventually $P > W$
- At this time efficiency E must drop because there will be processors doing no work
- Thus, problem size must grow at least by $\theta(P)$ for the problem to scale
 - $\theta(P)$ is the lower bound on the isoefficiency function
 - $\theta(P)$ is the isoefficiency function of an ideally scalable system

Degree of concurrency $C(W)$

- Lower bound of $\theta(P)$ for some algorithm is imposed by the algorithm's *degree of concurrency*
- If $\theta(P)$ is an algorithm's degree of concurrency, at most $\theta(P)$ processors can be used to solve the problem
- Example: Gaussian elimination has $\theta(n^3)$ amount of computation, but ...
 - n variables must be eliminated one after the other (sequentially)
 - n^2 work per variable
 - thus at most n^2 processors can ever be effectively be used at a time.

Degree of concurrency, cont.

- If $W = \theta(n^3)$ for this problem, degree of concurrency is $\theta(W^{2/3})$
- Given a problem of size W , at most $\theta(W^{2/3})$ processors can be used
 - For P processors, need $\theta(P^{3/2})$ work ($W^{2/3} = P$)
 - Thus, because of concurrency, isoefficiency function for this operation is $\theta(P^{3/2})$
- If algorithm's degree of concurrency is $< \theta(W)$, then
 - isoefficiency function due to concurrency is worse than $\theta(P)$
 - In these cases, isoefficiency function is the max of the isoefficiency functions due to concurrency, communication, and other overheads

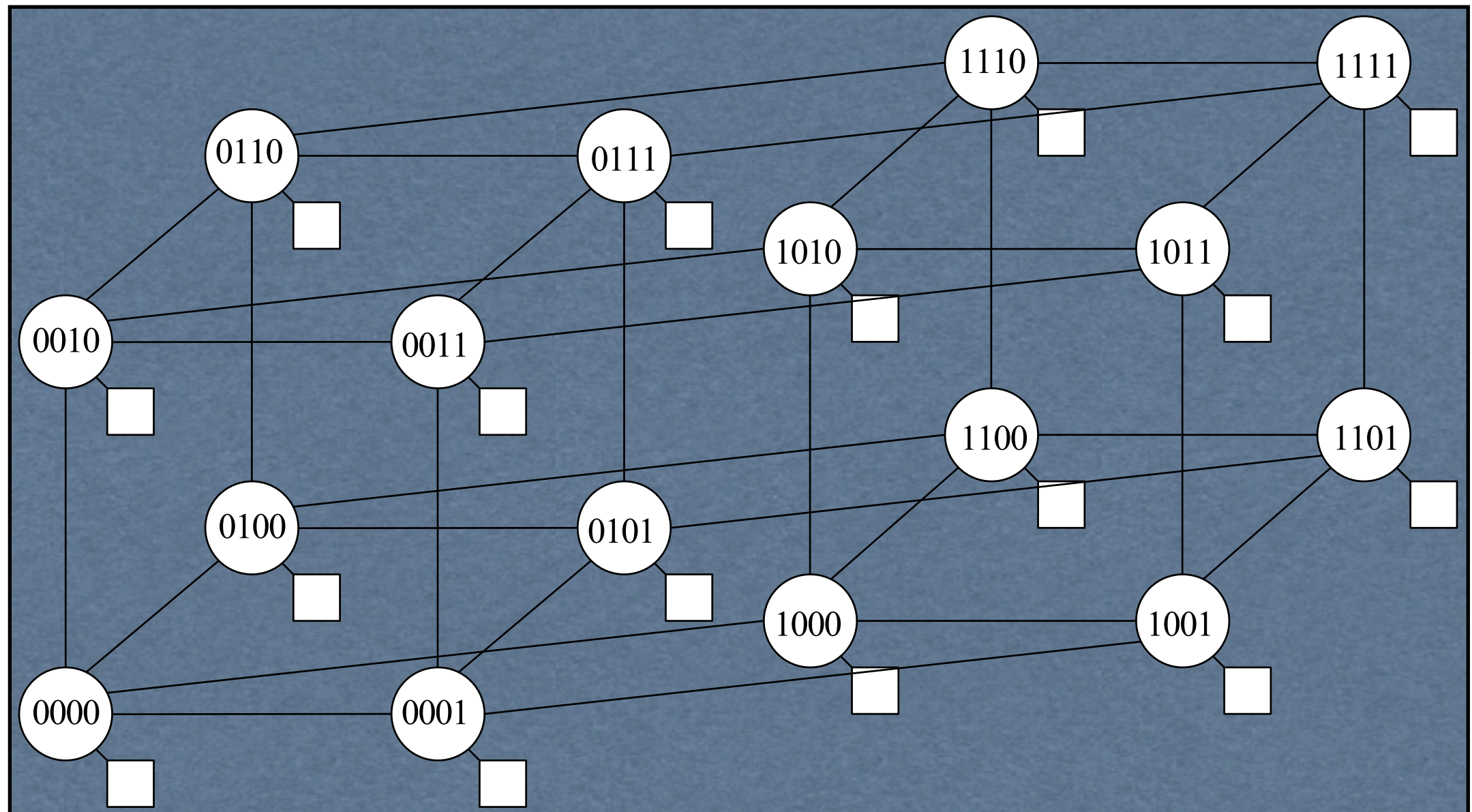
Hypercubes: short aside

- Since hypercubes are mentioned in the paper, let's talk about them for a few minutes.
- Hypercubes were first developed as part of the Cosmic Cube project at CalTech (Seitz and Fox). Commercial version came out as the Intel iPSC, with Cleve Moler as one of the designers.
- Cleve Moler went on to found Matlab, Jeff Fox now at IU CS, Seitz won 2011 IEEE Computer Society Seymour Cray Computer Engineering Award
- The original *Cosmic Cube* was a plot device used in Marvel Comics

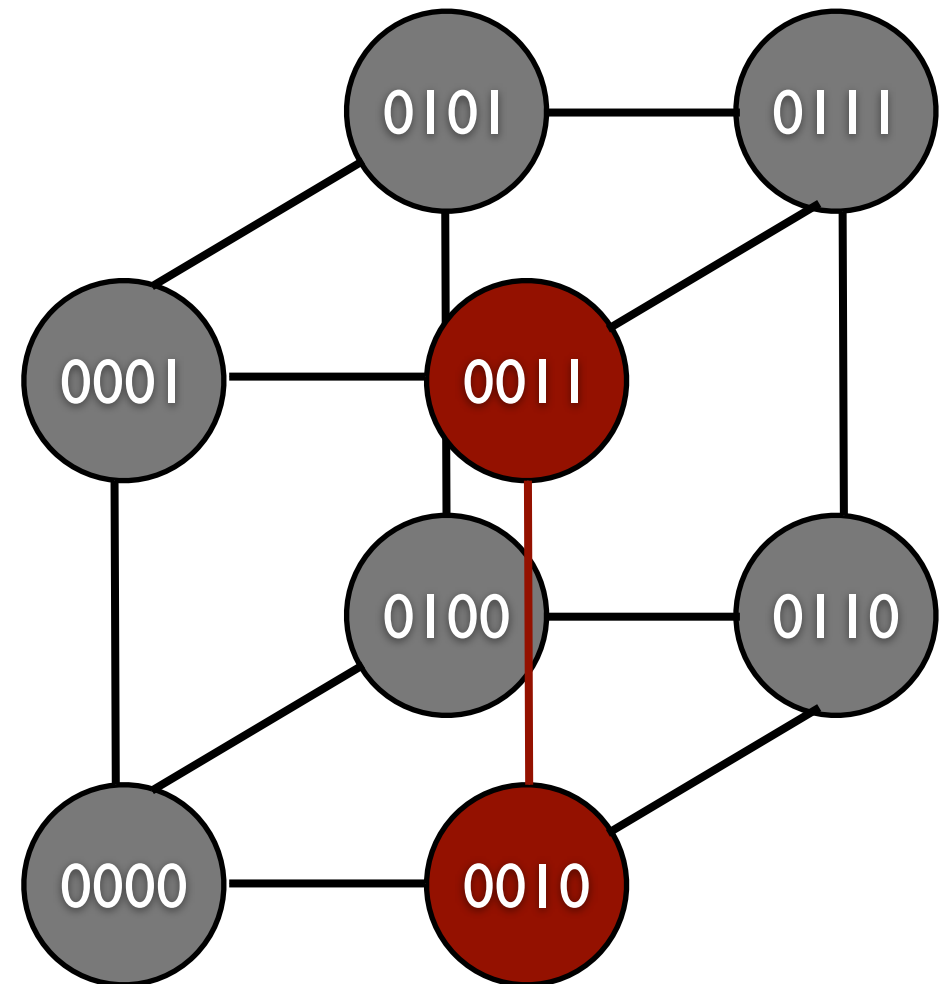
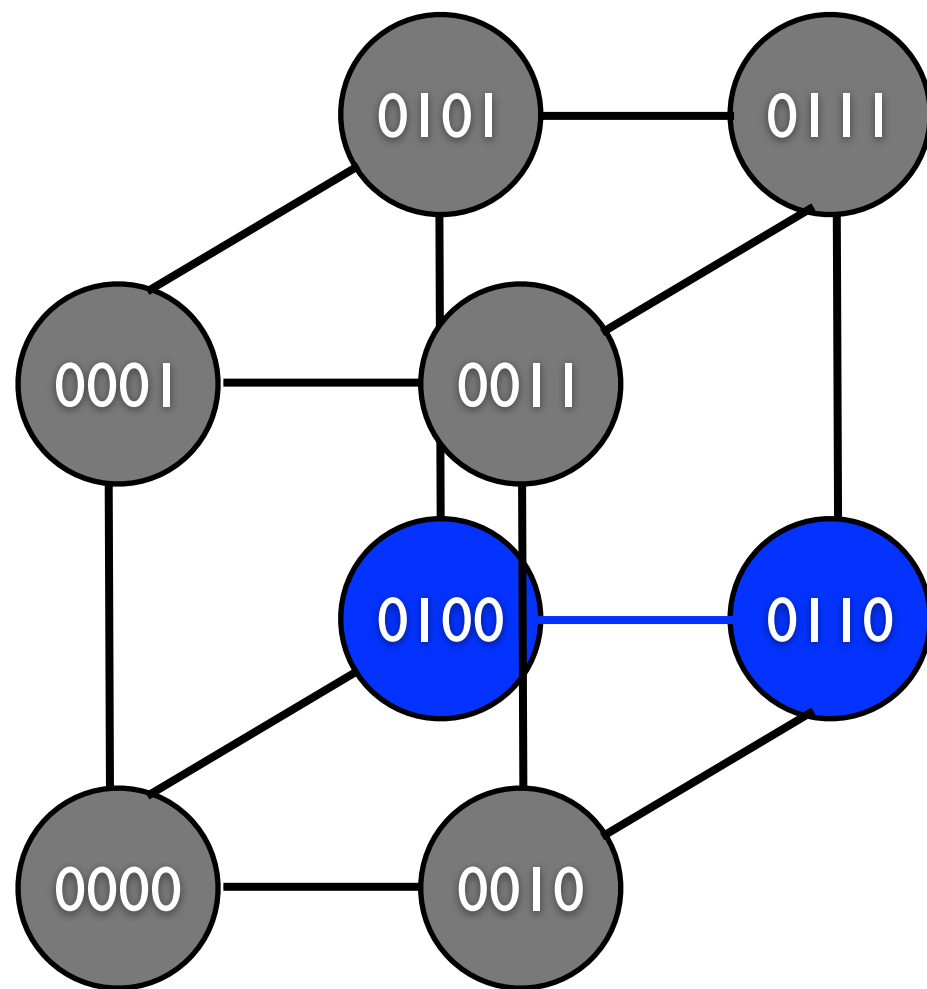
Hypercube

- Direct topology (one switch node/processor)
- $2 \times 2 \times \dots \times 2$ mesh
- Number of nodes a power of 2, denoted k
- Node addresses $0, 1, \dots, 2^k - 1$
- Node i connected to k nodes whose addresses differ from i in exactly one bit position

Hypercube Addressing



Hypercube labeling

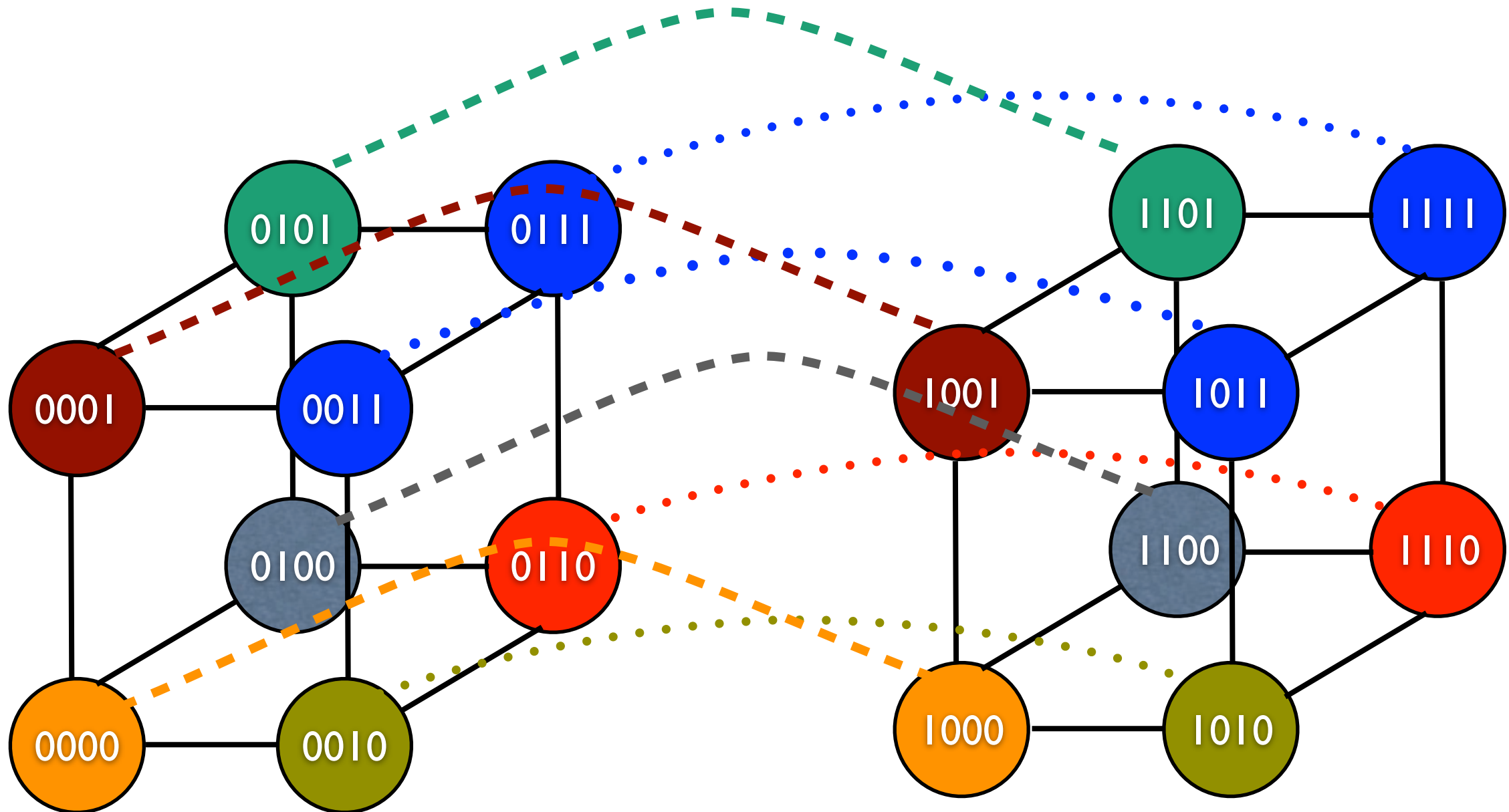


Pairs of adjacent nodes differ by 1 bit in their label --
result of gray code numbering

Hypercube labeling

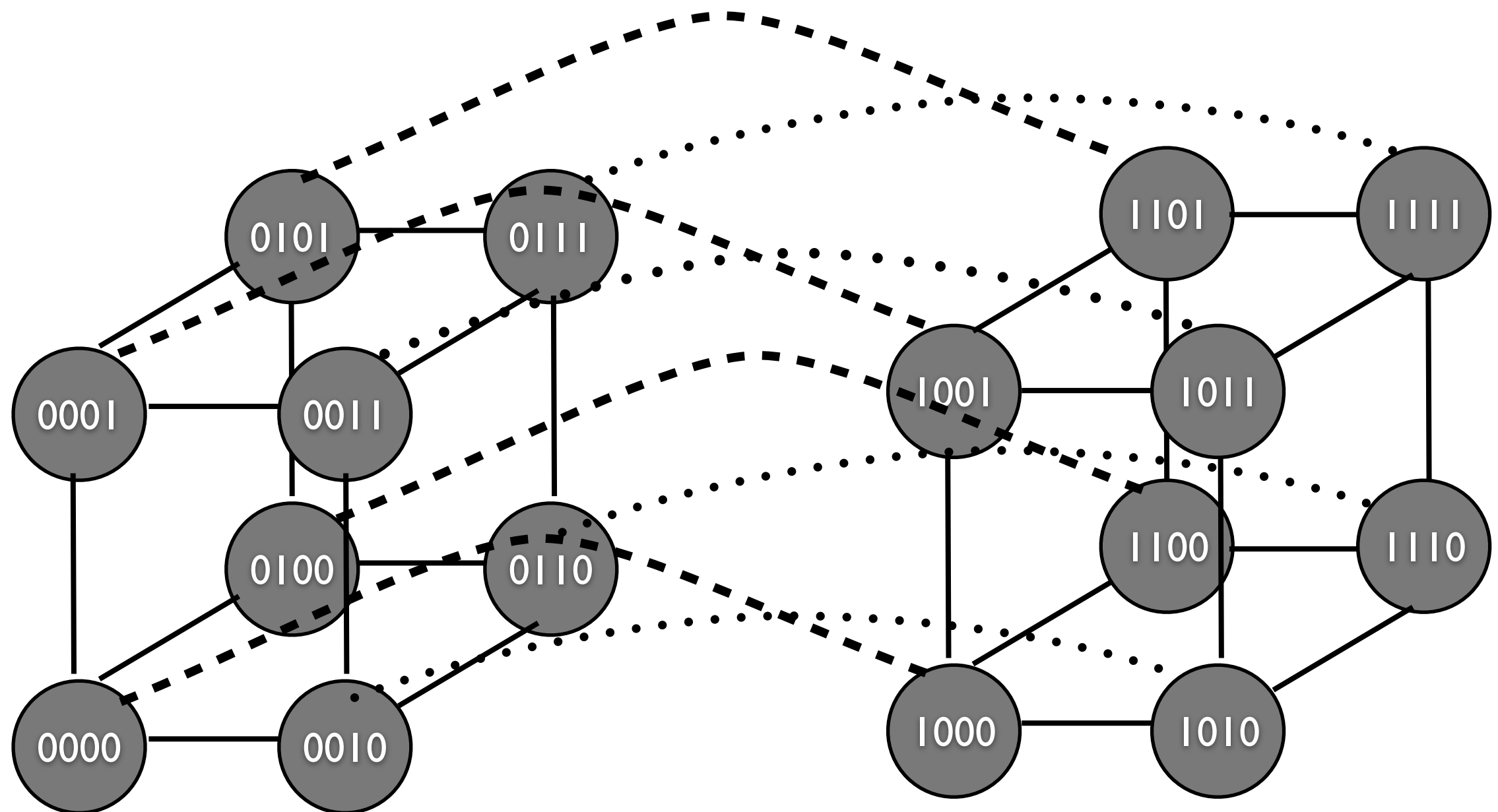
A large hypercube is made up of smaller hypercubes.

1. Add 1 bit to labels
2. Make bit 1 for one small hypercube, 0 for the other
3. Add edge to nodes whose labels differ by one



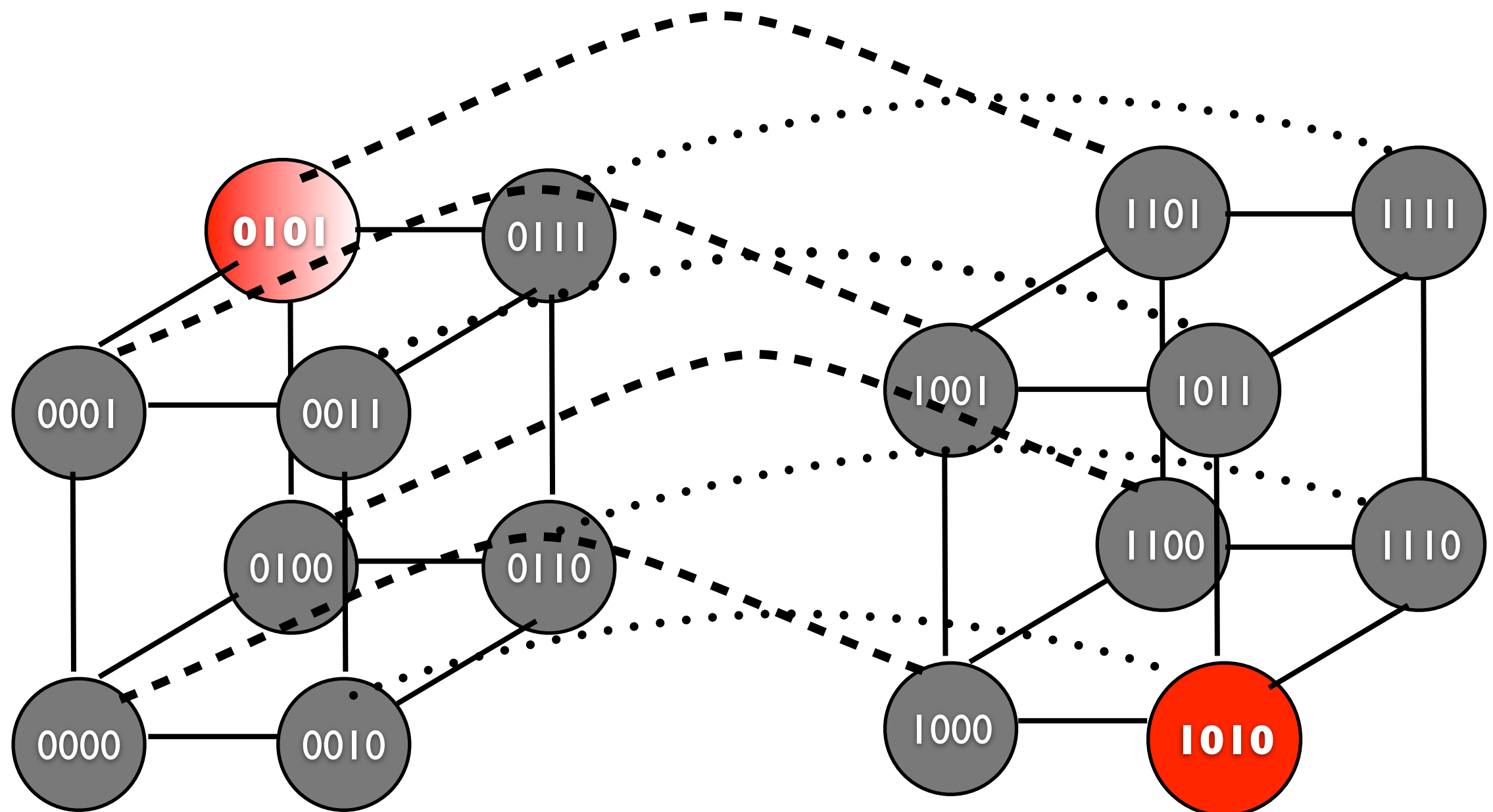
Labeling leads to routing

Given a source a destination label, always move one bit closer to the destination label with each hop.



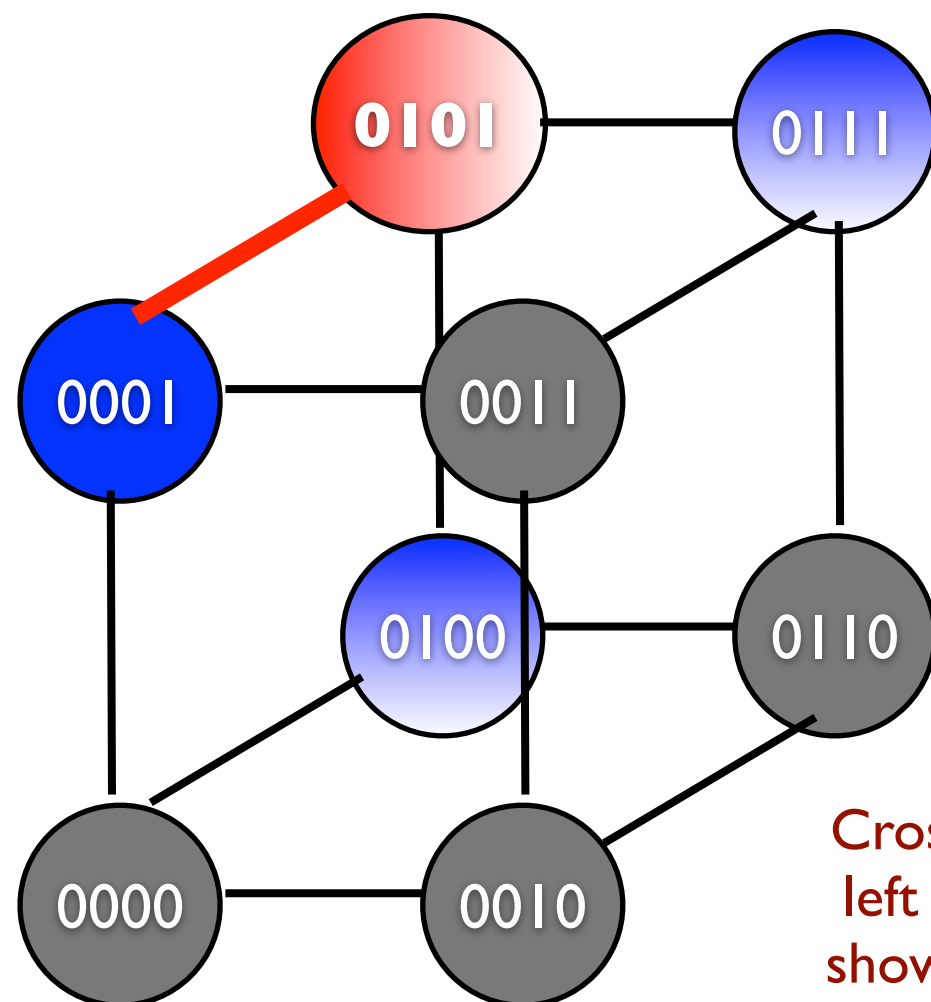
Labeling leads to routing

Go from 0101 to 1010, want to change source 0's to 1's, and 1's to 0's

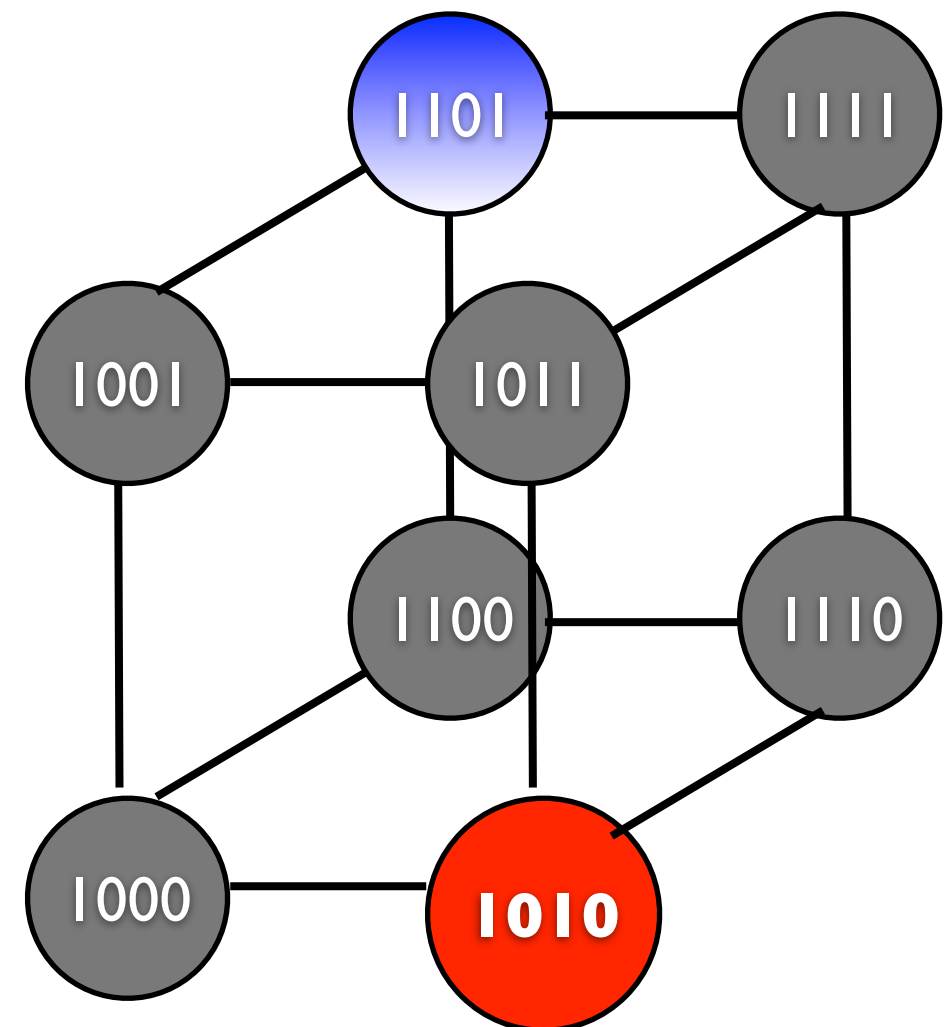


Labeling leads to routing

Go from 0101 to 0001, on the way to 1010. Note that since every bit needs to change, and every bit link changes one bit, we had four choices. In general, B choices, where B is the number of bits to change.

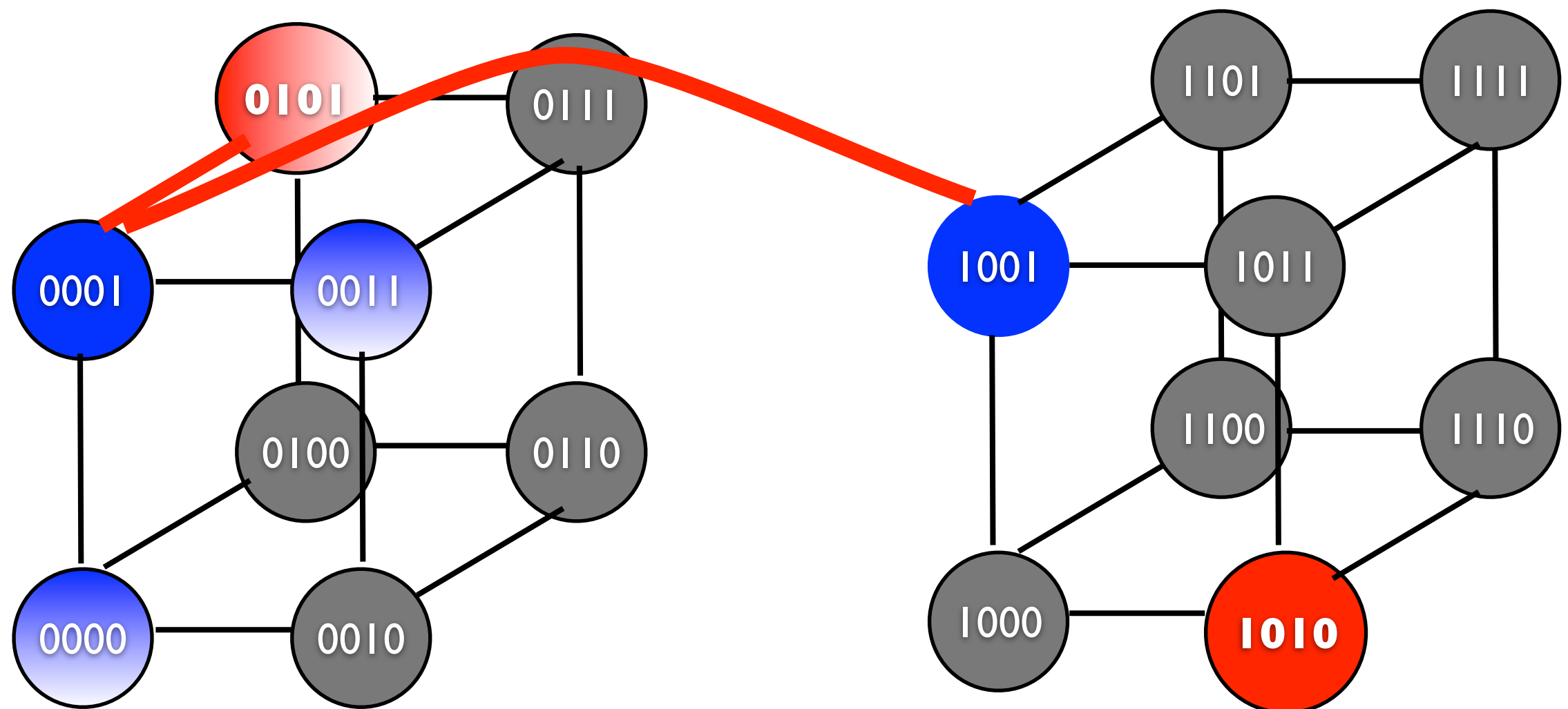


Cross links from
left to right not
shown for clarity



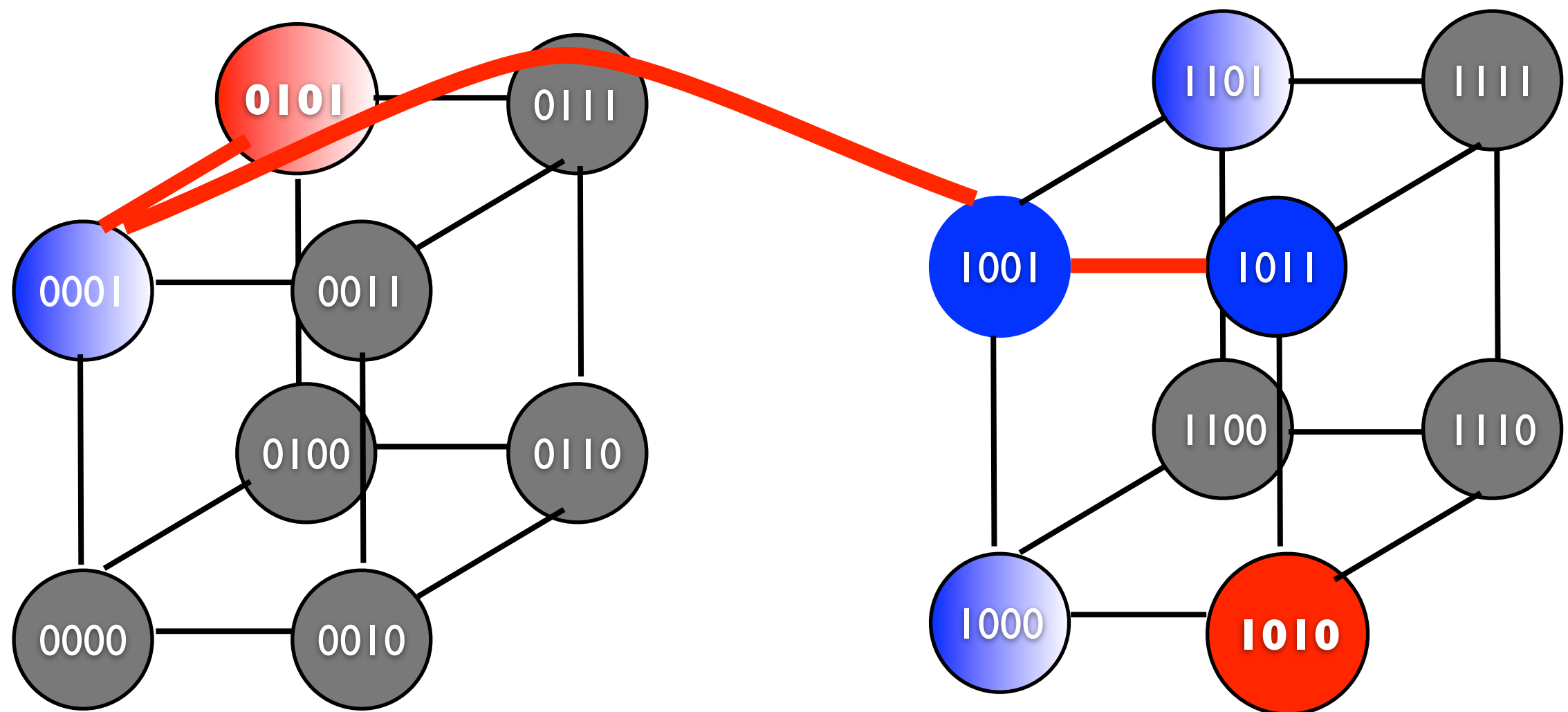
Labeling leads to routing

At **0001**, on route from 0101 to 1010. Three bits differ, three choices, pick one (**1001**)



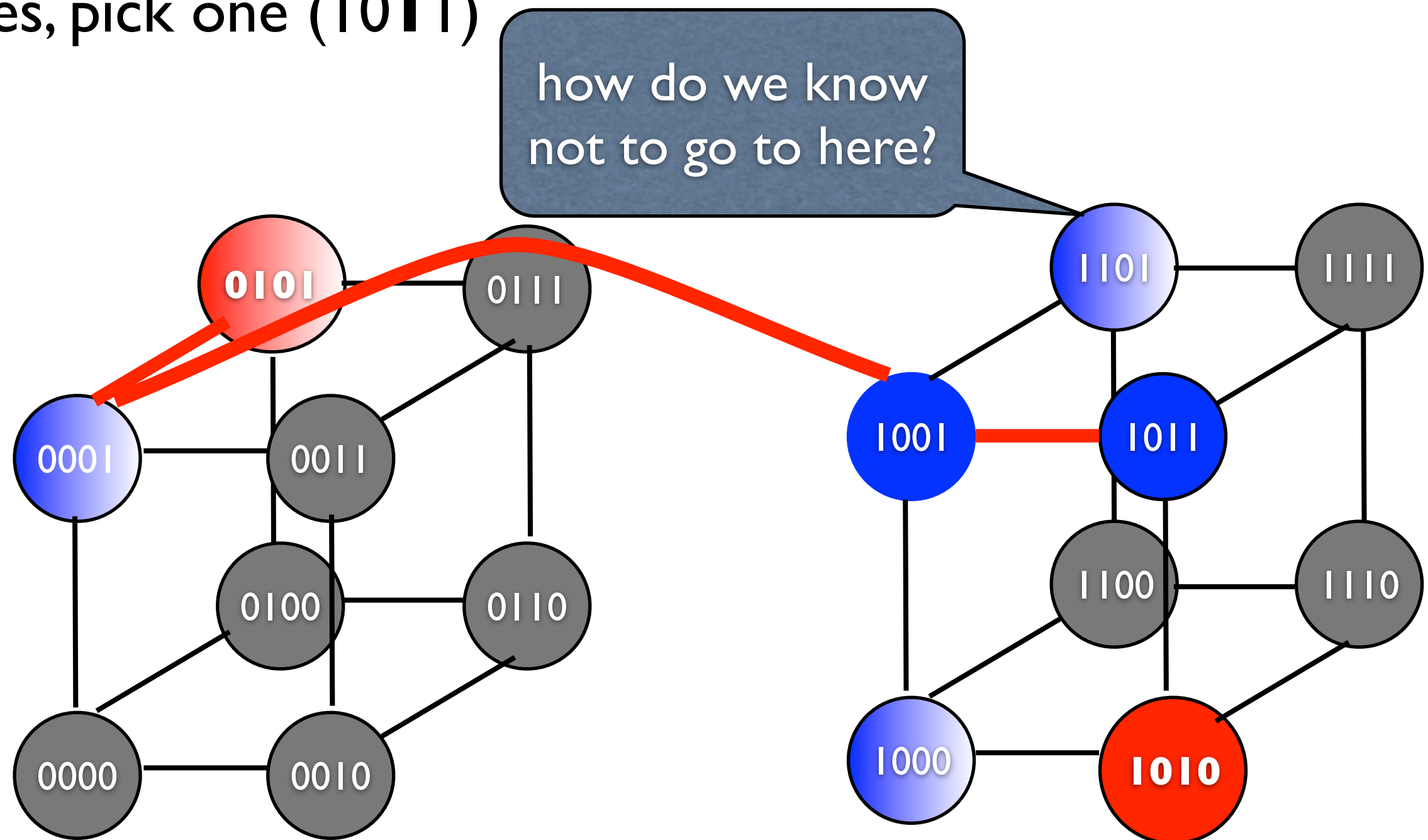
Labeling leads to routing

At **1001**, on route from 0101 to 1010. Two bits differ, two choices, pick one (1011)



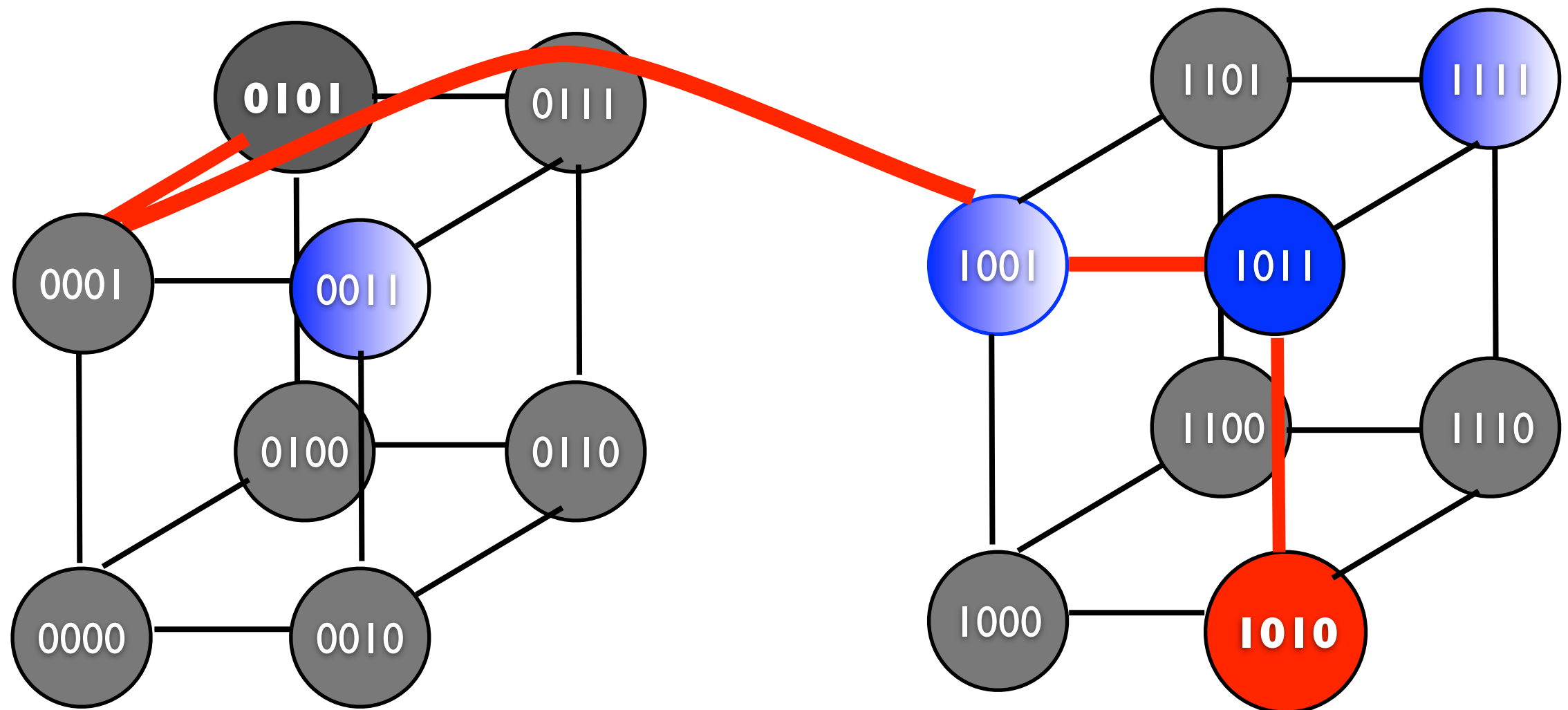
Labeling leads to routing

At **1001**, on route from 0101 to 1010. Two bits differ, two choices, pick one (1011)



Labeling leads to routing

At **1011**, on route from 0101 to 1010. One bit differs, only one choice, pick one (**1010**)



A slight aside

- Per-node performance is important.
- Cache and prefetch effects are an important way to gather per-node performance.

Loop unrolling

```
do i = 1, n  
  a[i] = b[i] + c[i]  
end do
```

- Processors typically have several add and store units
- Without additional hardware, processors cannot move operations around conditional branches
- Thus, in this loop, only one add and one store can be started at a time

```
do i = 1, n, 3  
  a[i] = b[i] + c[i]  
  a[i+1] = b[i+1] + c[i+1]  
  a[i+2] = b[i+2] + c[i+2]  
end do
```

```
do i = n-((n-1) mod 3), n  
  a[i] = b[i] + c[i]  
end do
```

- In the right hand version, three floating point adds can be issued at once
- Blue fixup code necessary because n not always divisible by the unroll factor

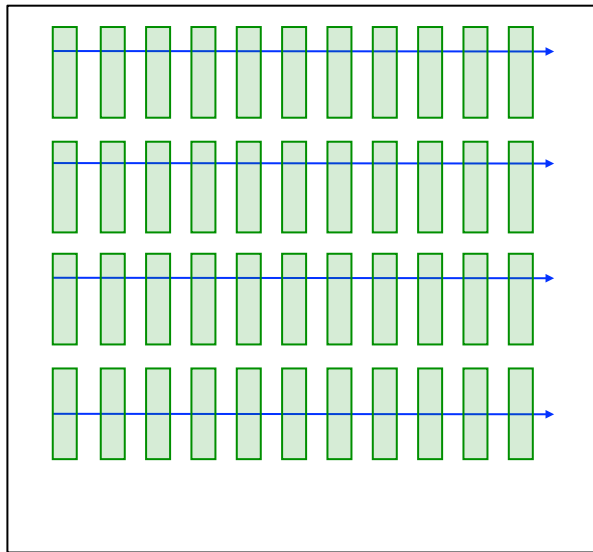
Cache optimizations

- Desire is to maximize locality and to exploit temporal and spatial locality via the cache.
- Accesses from cache are tens to hundreds of times faster than accesses from memory, and typically the same speed or 2 times slower than accesses from registers
- The more regular the code, the more compilers can do in terms of cache optimizations

Cache optimizations

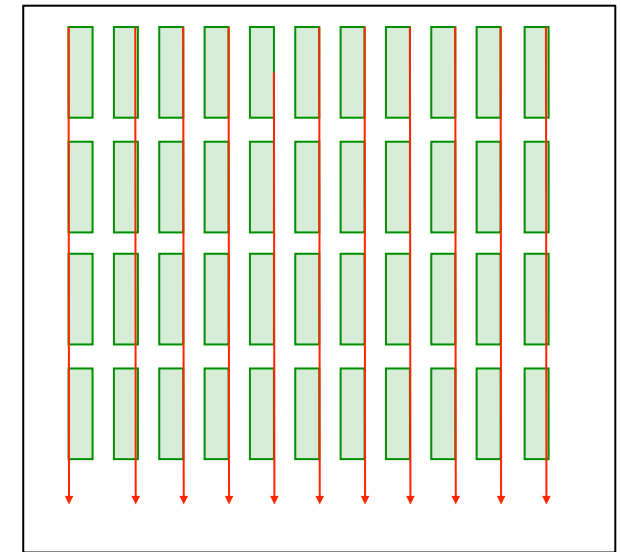
- Loop interchange attempts to make accesses coincide with the way data is laid out in memory

```
do i = 1, n  
  do j = 1, n  
    a[i,j] = ...  
  end do  
end do
```



A[col,row], array stored in
column major order (row in
C)

```
do j = 1, n  
  do i = 1, n  
    a[i,j] = ...  
  end do  
end do
```

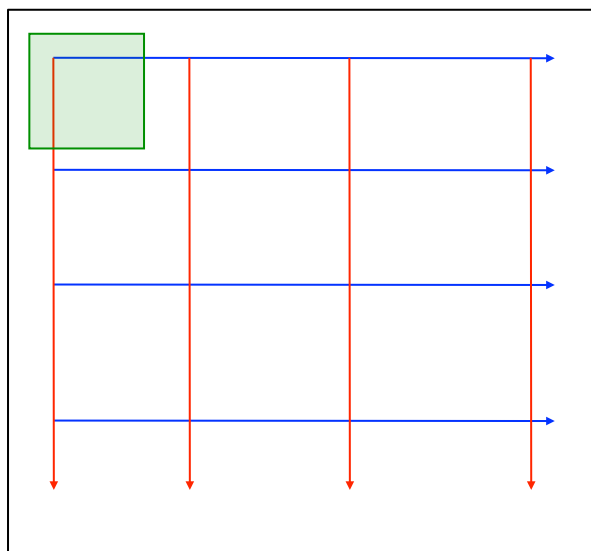



Interchange the loops, and data within a
text box is reused.

Loop tiling (page 347, Bacon, Eggers paper on the class web page)

```
Do i = 1, n
  Do j = 1, n
    a[i,j] = b[j,i]
  end do
end do
```

a,b[col,row], array stored
in row major order




What fits in
cache

```
do TI = 1, n, 64
  do TJ = 1, n, 64
    do i = TI, min(TI+63,n)
      do j = TJ, min(TJ+63, n)
        a[i,j] = b[j,i]
      end do
    ...
  end do
```

b[i,j] accesses

a[i,j] accesses

This pattern is present in matrix
multiply, transpose, etc.

Cache optimizations

- Loop fusion increases temporal locality

```
do l = 1, n
```

```
  c[i] = a[i]
```

```
end do
```

```
do l = 1, n
```

```
  b[i] = a[i]
```

```
end do
```

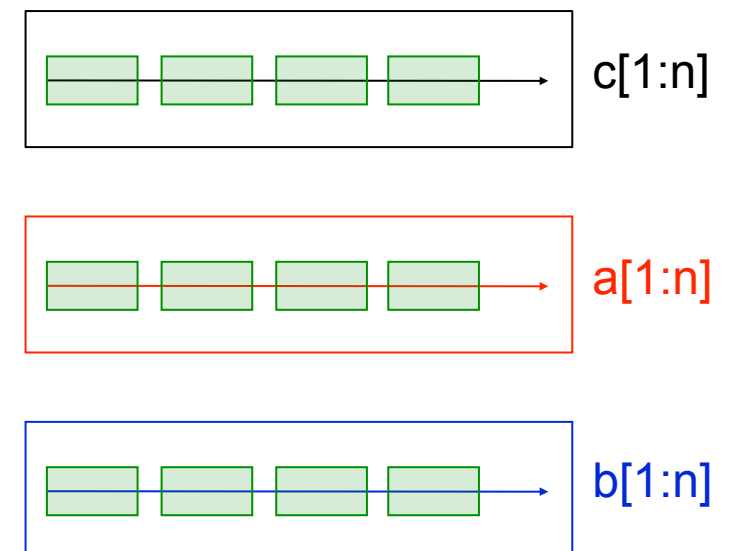
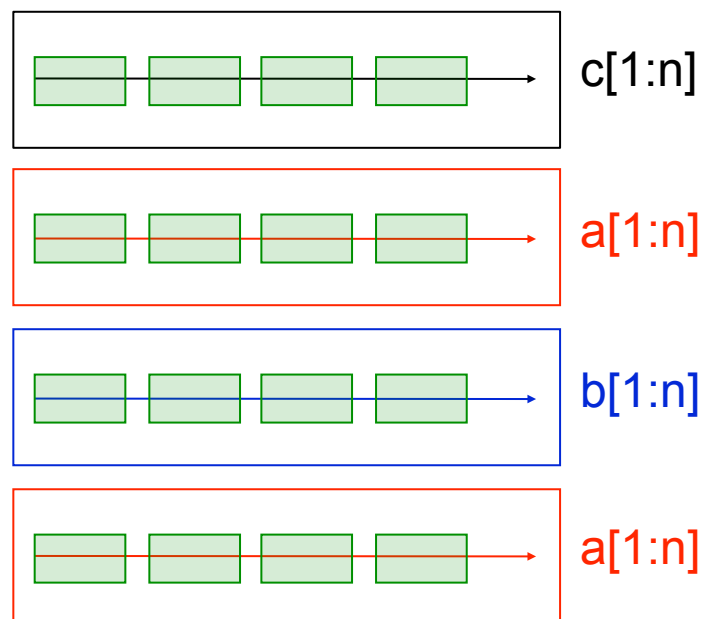
If n is large enough that all of $c[1:n]$, $a[1:n]$ and $b[1:n]$ won't fit in cache, some of $a[1:n]$ fetched in the first loop will be evicted by the time we need it in the second loop, and need to be refetched.

```
do l = 1, n
```

```
  c[i] = a[i]
```

```
  b[i] = a[i]
```

```
end do
```



Cache optimizations

- Loop fusion, like all transformations that change the access order of storage, need to be checked for legality

```
do i = 1, n
  c[i] = a[i]
end do
do i = 1, n
  b[i] = c[i+1]
end do
```

For example, iteration 4 of the second loop reads $c[5]$, which is not written until iteration 5 of the first loop. The fused loop (on the right) reads stale values for c .

By shifting the iteration spaces we get a correct execution, but the transformation gets more complicated.

```
do i = 1, n
  c[i] = a[i]
  b[i] = c[i+1]
end do
```

```
c[1] = a[1]
do i = 2, n
  c[i] = a[i]
  b[i-1] = c[i]
end do
b[n] = c[n+1]
```

Cache optimizations

- Prefetching attempts to get values in cache before they are used

```
do i = 1, n
  do j = 1, n
    a[i,j] = b[i,j]
  end do
end do

load a[1,1], b[1,1]
...
do TI = 1, n, 64
  do TJ = 1, n, 64
    load a[TI,TJ], b[TI,TJ]
    do i = TI, min(TI+63,n)
      do j = TJ, min(TJ+63, n)
        a[i,j] = b[j,i]
      end do
    end do
    ...
  end do
```

Problems with prefetching

- Need to fetch early enough to make a difference
- What about page faults?
- What if the value is changed between prefetch and use?
- What if we **evict** a needed value from the cache?

Cache optimizations

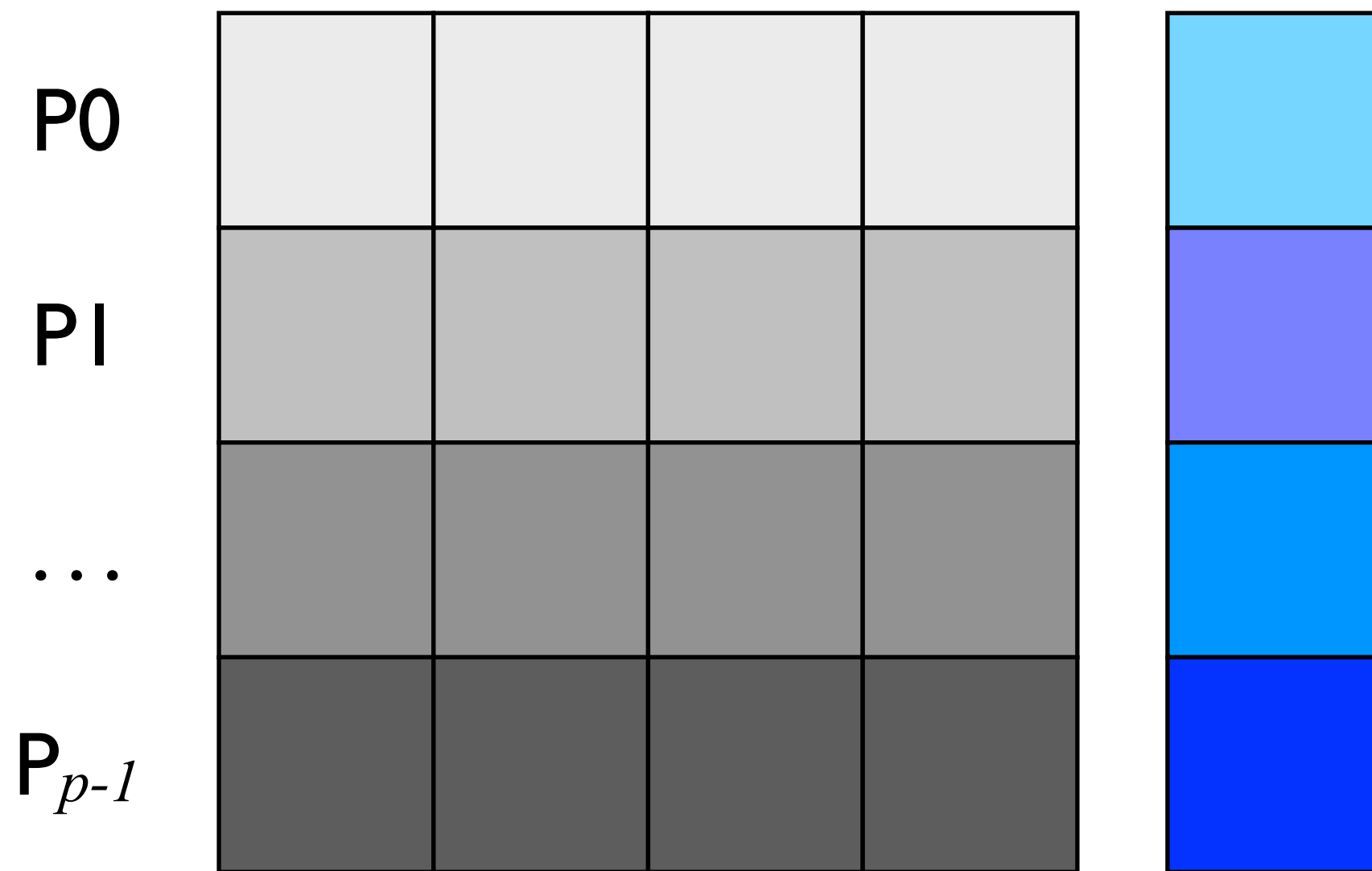
Array padding eliminates or reduces false sharing across cache lines.

See Bacon and Eggers for details.

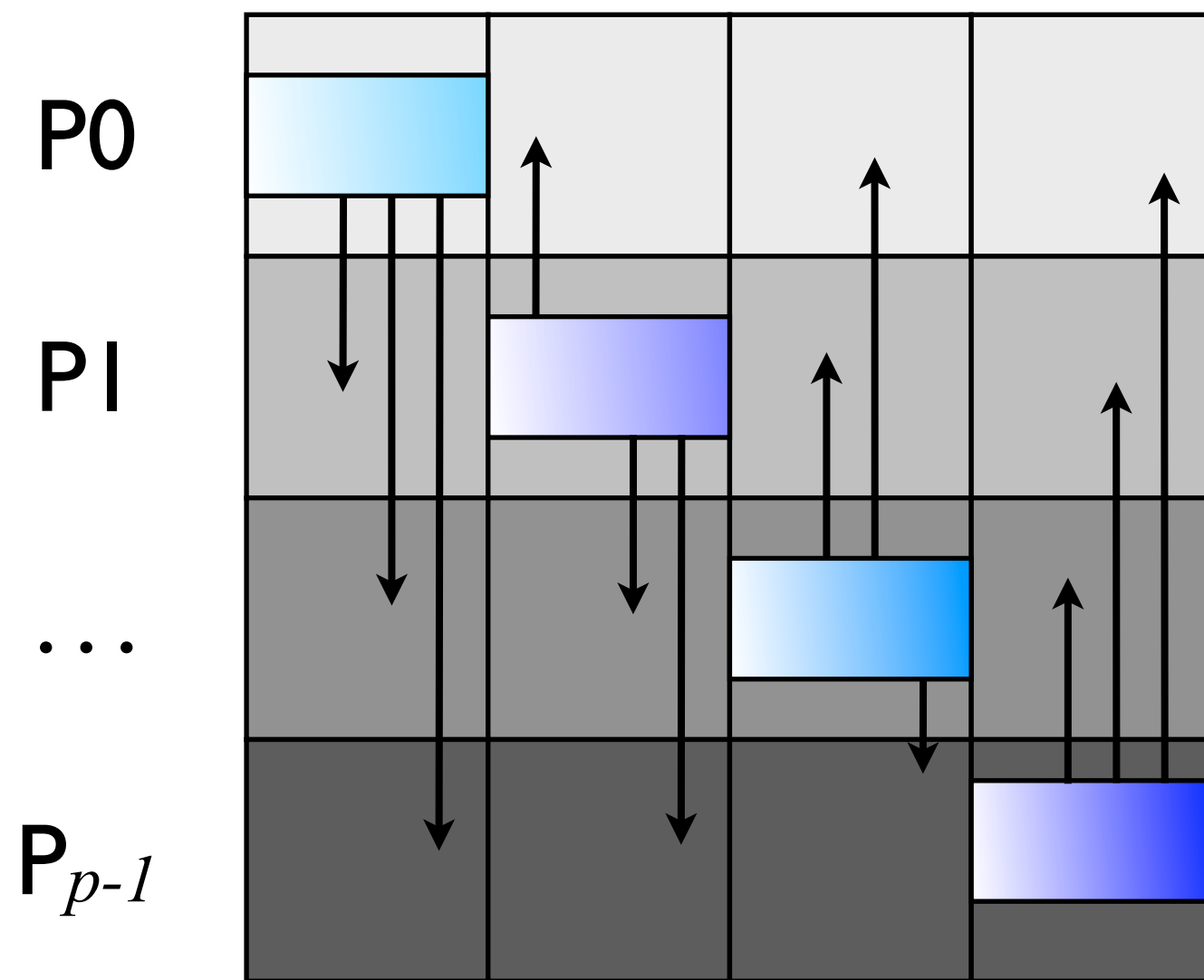
Comparing matrix vector algorithms - sequential alg.

- Consider matrix vector multiply, i.e. an $n \times n$ matrix times an $n \times 1$ matrix
- Number of basic operations (W) is n^2 , with t_c the time for a single floating multiply-add
- Sequential time is $n^2 t_c$, i.e. $T_1 = n^2 t_c$

With striped, data starts out like this, e.g.,
from reading matrix and vector from disks



Every process sends its n/P elements of the vector to every other process



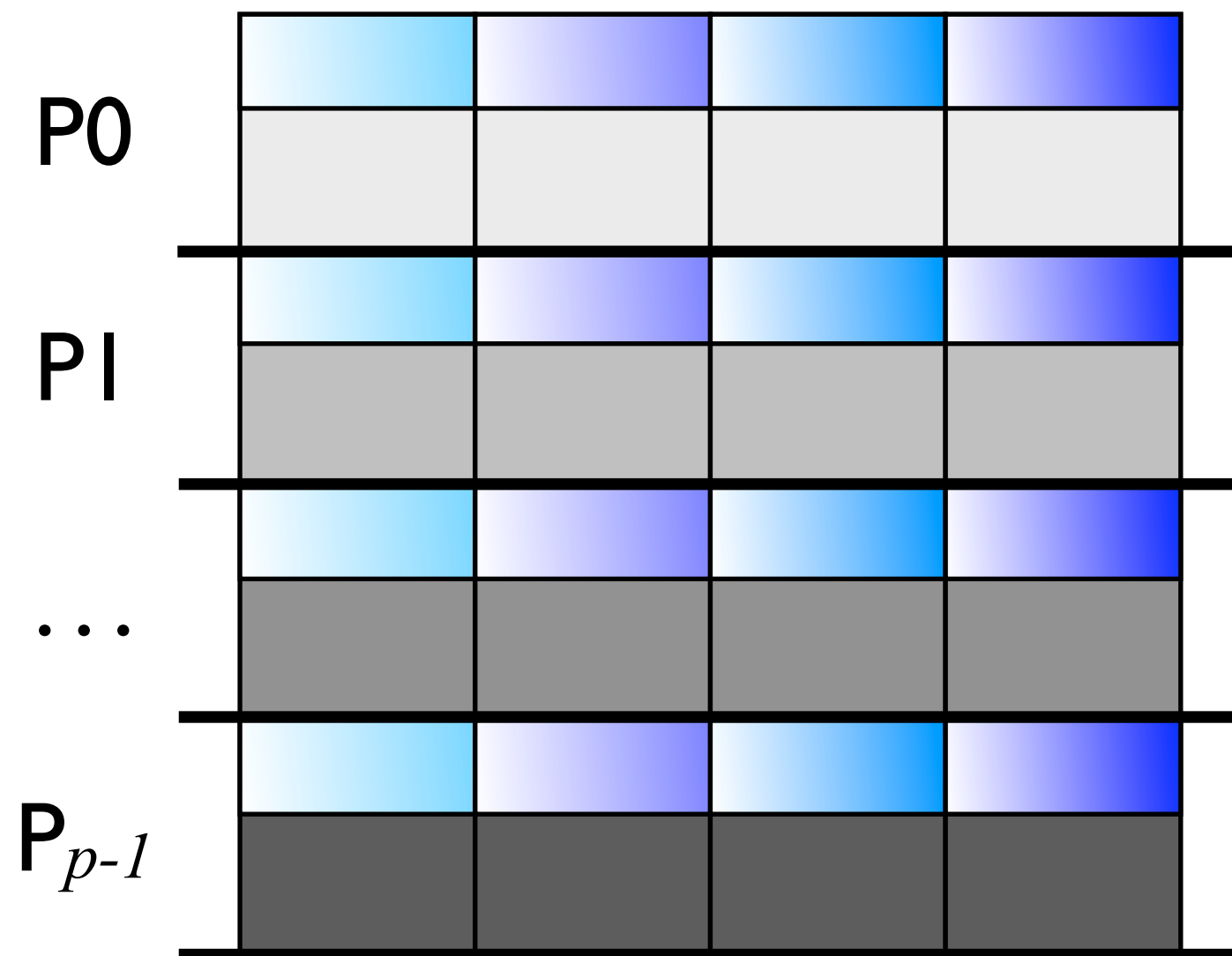
Do this with an all-to-all broadcast

$$t_s \log P + t_w n(P-1)/P$$

startup
time (t_s is
startup
time of the
network

comm. time
-- t_w is the
time to send
one word

After communication every process has a copy of the vector



Row-striped parallel alg.

- n/p matrix rows and vector elements to each processor
- Costs:
 - all-to-all broadcast of vector elements so that each processor has a copy:

$$t_s \log P + t_w n(P-1)/P$$

or simply

$$t_s \log P + t_w n$$

as P grows large (where t_s is startup time, t_w is per-word transfer time)

Row-striped parallel alg.

- n/p matrix rows and vector elements to each processor
 - Each node does $t_c n^2/p$ work multiplying n/p rows times the vector
 - $T_P = t_c n^2/P + t_s \log P + t_w n$
 - Using the relation $T_O = P T_P - T_1$, we get $T_O = t_s P \log P + t_w n P$

$$T_O = P(t_c n^2/P + t_s \log P + t_w n) - n^2 t_c$$

$$= t_c n^2 + t_s P \log P + t_w P n - n^2 t_c$$

$$= t_s P \log P + t_w n P$$

Isoefficiency relationship

- $T_O = t_s P \log P + t_w n P$
- Balance the first term of T_O by rewriting $W = K T_O$ using only first term of $T_O = t_s P \log P + t_w n P$ to get $W = K t_s P \log P$
- Balancing the second term of T_O (due to per-word transfer time) against the problem size W and in terms of P we get

$$n^2 = K t_w n P$$

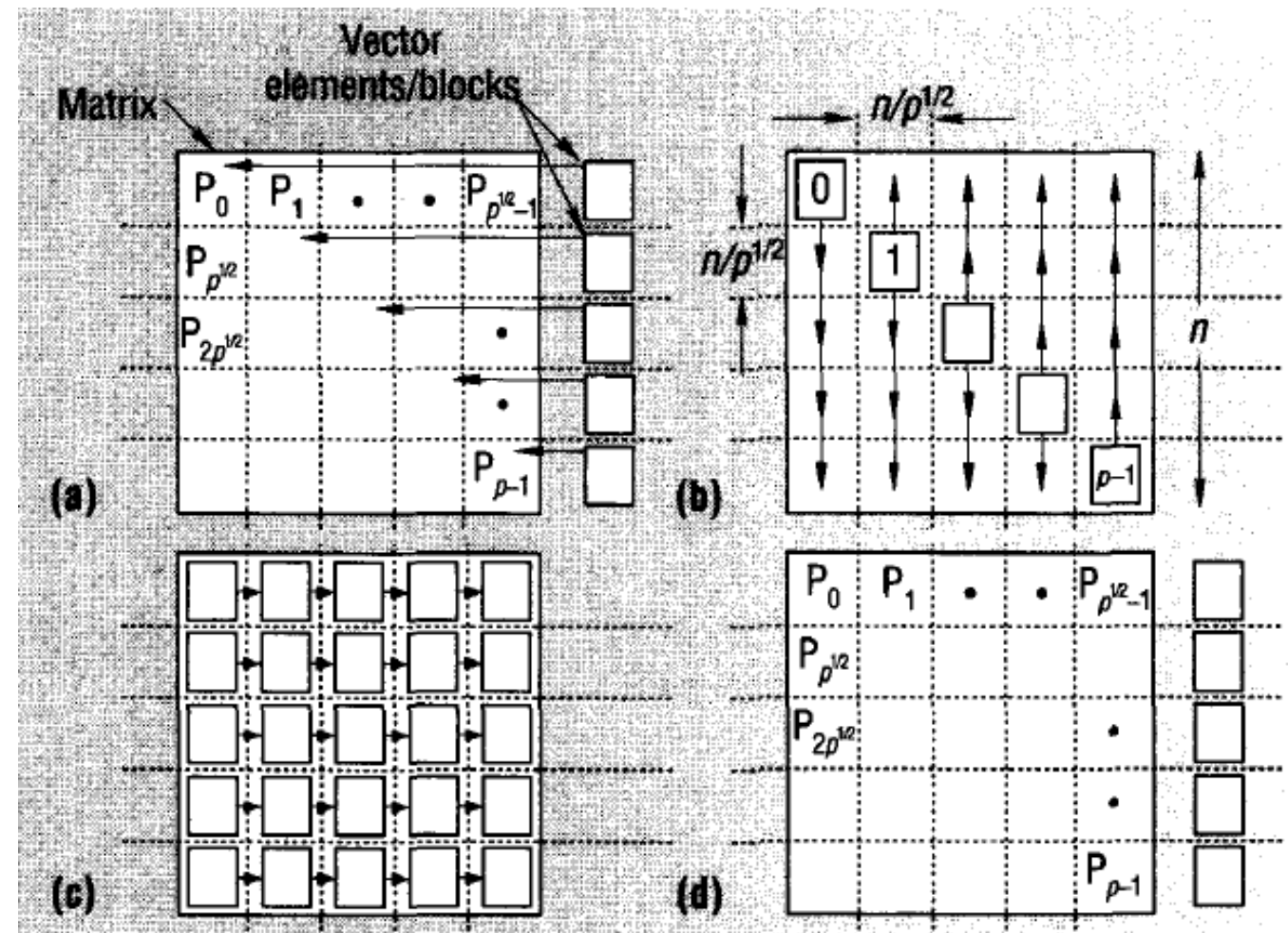
$$n = K t_w P \text{ (solve for } n \text{ in terms of } K \text{ and } t_w \text{ (constants) and } P\text{)}$$

$$W = n^2 = K^2 t_w^2 P^2$$

- To maintain efficiency, work must increase proportional to P^2

Checkerboard partitioning - data is originally in the last processor of each column

- Divide data into $n/\sqrt{p} \times n/\sqrt{p}$ squares and place on the last column of processes
- Each process w/data sends it to the diagonal of its row (a)
- Column-wise one-to-all broadcast of n/\sqrt{p} elements (b)
- Each processor performs n^2/p multiplications, and locally adds n/\sqrt{p} sets of products. (c)
- n/\sqrt{p} partial sums to be accumulated along each row (c)



State at end of computation (d)

Checkerboard partitioning analysis

1. Divide data into $n/\sqrt{p} \times n/\sqrt{p}$ squares, send along rows
2. Column-wise one-to-all broadcast of n/\sqrt{p} elements
takes $t_s + t_w(n/\sqrt{p}) \log \sqrt{p}$ time on a hypercube with store-and-forward routing, $2 (t_s + t_w(n/\sqrt{p}) \log \sqrt{p})$ time total for 1 & 2
3. Each processor performs n^2/p multiplications, and locally adds n/\sqrt{p} sets of products.
takes $t_c n^2/p$ time
4. n/\sqrt{p} partial sums to be accumulated along each row
also takes $t_s + t_w(n/\sqrt{p}) \log \sqrt{p}$ time on a hypercube with store-and-forward routing using a reduction
5. total parallel time is
$$T_P = t_c(n^2/p) + t_s + 2 t_s \log \sqrt{p} + 3 t_w (n/\sqrt{p}) \log \sqrt{p}$$

Simplify

- Can approximate

$$T_P = t_c(n^2/p) + t_s + 2 t_s \log \sqrt{p} + 3 t_w (n/\sqrt{p}) \log \sqrt{p}$$

- with (substituting $(\log p)/2$ for $\log \sqrt{p}$), ignoring non-p terms

$$T_P = t_c(n^2/P) + t_s \log p + (3/2) t_w (n/\sqrt{p}) \log p$$

serial
work

- will use this expression to find isoefficiency, in particular

$$T_O = \cancel{t_c(n^2/p)} + t_s p \log p + (3/2) t_w (n/\sqrt{p}) \log p$$

- and thus $T_O = t_s p \log p + (3/2) t_w (n/\sqrt{p}) \log p$

Simplify and analyze

$$T_O = t_s p \log p + (3/2) t_w n \sqrt{p} \log p$$

- Solve for isoefficiency resulting from the t_w term
Equate each term of T_O with the problem size W in terms of P and constants

$$n^2 t_c = K (3/2) t_w n \sqrt{p} \log p$$

$$n = K (3/2) (t_w/t_c) \sqrt{p} \log p$$

$$W = n^2 = K (9/4) (t_w^2/t_c^2) p \log^2 p$$

constants for a
given problem and
machine

- The isoefficiency due to t_w is $\theta(p \log^2 p)$
- This is also overall isoefficiency, since it dominates the $\theta(p \log p)$ term involving t_s

What we can conclude

- For the striped model $W = n^2 = K^2 t_w^2 P^2$, and to maintain efficiency, work must increase proportional to P^2
- For the checkerboard model, $\theta(p \log^2 p)$ and $p \log^2 p < P^2$
- Therefore, the checkerboard model will scale better than the striped model
- The fundamental reason for this is that the communication is over a smaller number of processors

Isoefficiency and concurrency

- Some algorithms with low overhead also have limited concurrency
- This has a negative effect on isoefficiency, as we will see from Dijkstra's all-pairs shortest-path algorithm
- One instance of Dijkstra's algorithm computes the shortest distance between a single node **s** and all other nodes

Edgar Dijkstra

- Dutch computer scientist, eventually worked to UT Austin, didn't particularly like computers, considered fairly cranky (but very smart and dedicated to teaching) by those who worked with him.

The job [of operating or using a computer] was actually beyond the electronic technology of the day, and, as a result, the question of how to get and keep the physical equipment more or less in working condition became in the early days the all-overriding concern. As a result, the topic became —primarily in the USA— prematurely known as "computer science" —which, actually is like referring to surgery as "knife science"— and it was firmly implanted in people's minds that computing science is about machines and their peripheral equipment. Quod non [Latin: "Which is not true"]

- “And I don’t need to waste my time with a computer just because I’m a computer scientist. [Medical researchers are not required to suffer from the diseases they investigate.]” EWD 1305

Edgar Dijkstra

I think anthropomorphism is worst of all. I have now seen programs "trying to do things", "wanting to do things", "believing things to be true", "knowing things" etc. Don't be so naive as to believe that this use of language is harmless. It invites the programmer to identify himself with the execution of the program and almost forces upon him the use of operational semantics.

Edgar Dijkstra

We could, for instance, begin with cleaning up our language by no longer calling a bug a bug but by calling it an error. It is much more honest because it squarely puts the blame where it belongs, viz. with the programmer who made the error. The animistic metaphor of the bug that maliciously sneaked in while the programmer was not looking is intellectually dishonest as it disguises that the error is the programmer's own creation... My next linguistically suggestion is more rigorous. It is to fight the "if-this-guy-wants-to-talk-to-that-guy" syndrome: never refer to parts of programs or pieces of equipment in anthropomorphic terminology...

EMD books, Dijkstra font

reasoning steps needed to keep the design under strict intellectual control. What is needed to achieve this goal, I can only describe as improving one's mathematical skills, where I use mathematics in the sense of "the art and science of effective reasoning". As a matter of fact, the challenges of designing high-quality programs and of designing high-quality proofs are very similar, so similar that I am no longer able to distinguish between the two: I see no meaningful difference between programming methodology and mathematical methodology

I came across a [comment on Reddit](#) by someone that had Dijkstra as a professor. Here's what it said:

I've always had horrible handwriting. When I was a computer science student I was in a class taught by Edsger Dijkstra. During the class he asked us to occasionally turn in our notes, because he wanted to see what we thought was important.

The final was an oral final and after going through a few questions to his satisfaction he said "You seem competent, but your handwriting is horrible..." The remaining 30 mins of my final exam by Dijkstra was me writing phrases repeatedly on a pad of paper while he said, 'no, you need to round the o's a bit more, the A is misformed, etc...'.

It was surreal. I'm sad he died.

From <https://joshldavis.com/2013/05/20/the-path-to-dijkstras-handwriting/>

Incidentally, there is a possibly apocryphal story that when Luca Cardelli used this font in a presentation at which Dijkstra was present, his comment to Cardelli was that he was impressed by how improved his handwriting was.

AaBbCcDdEeFfGgHhIi
JjKkLlMmNnOoPpQqRr
SsTtUuVvWwXxYyZz
1234567890

Dijkstra's algorithm

// d_i is the distance from d_s to d_i
// V is the set of N vertices

1. procedure sequential_dijkstra
2. $d_s = 0$
3. $d_i = \infty, i \neq s, i \in V$
4. $T = V$
5. for $i = 0$ to $N-1$
6. find $v_m \in T$ with minimum d_m
7. for each edge (v_m, v_t) with $v_t \in T$
8. if $(d_t > d_m + \text{length}((v_m, v_t)))$ then
9. $d_t = d_m + \text{length}((v_m, v_t))$
10. $T = T - v_m$

To find the shortest path from a vertex s to all other vertices

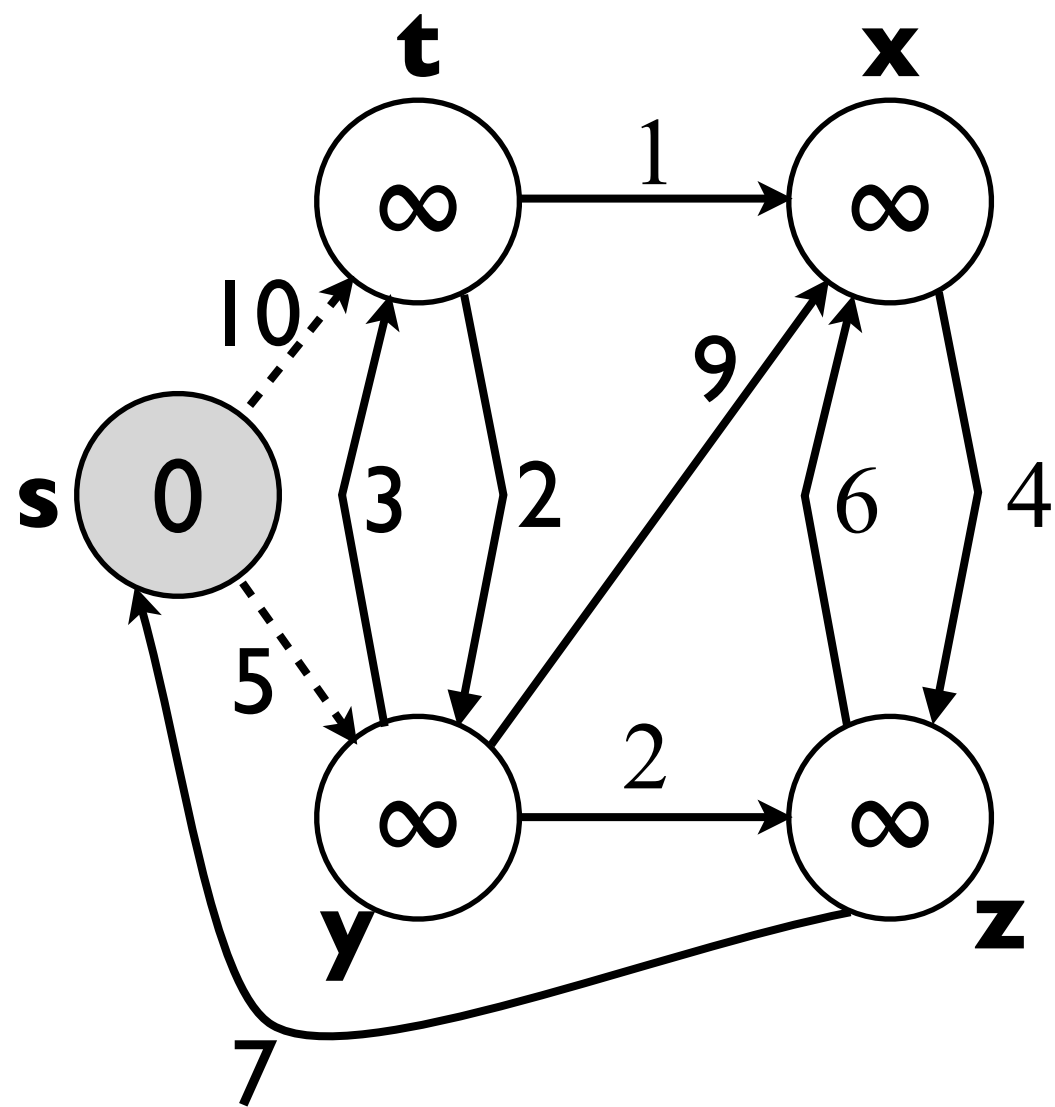
At each step

pick the node to be processed (a member of T)
 v_m that is closest to v_m (this is v_m on the first iteration)

for every other node v_t that is to be processed
see if there is a edge from v_m to v_t that leads to a shorter distance from s to v_t

remove v_m from the set of unprocessed nodes

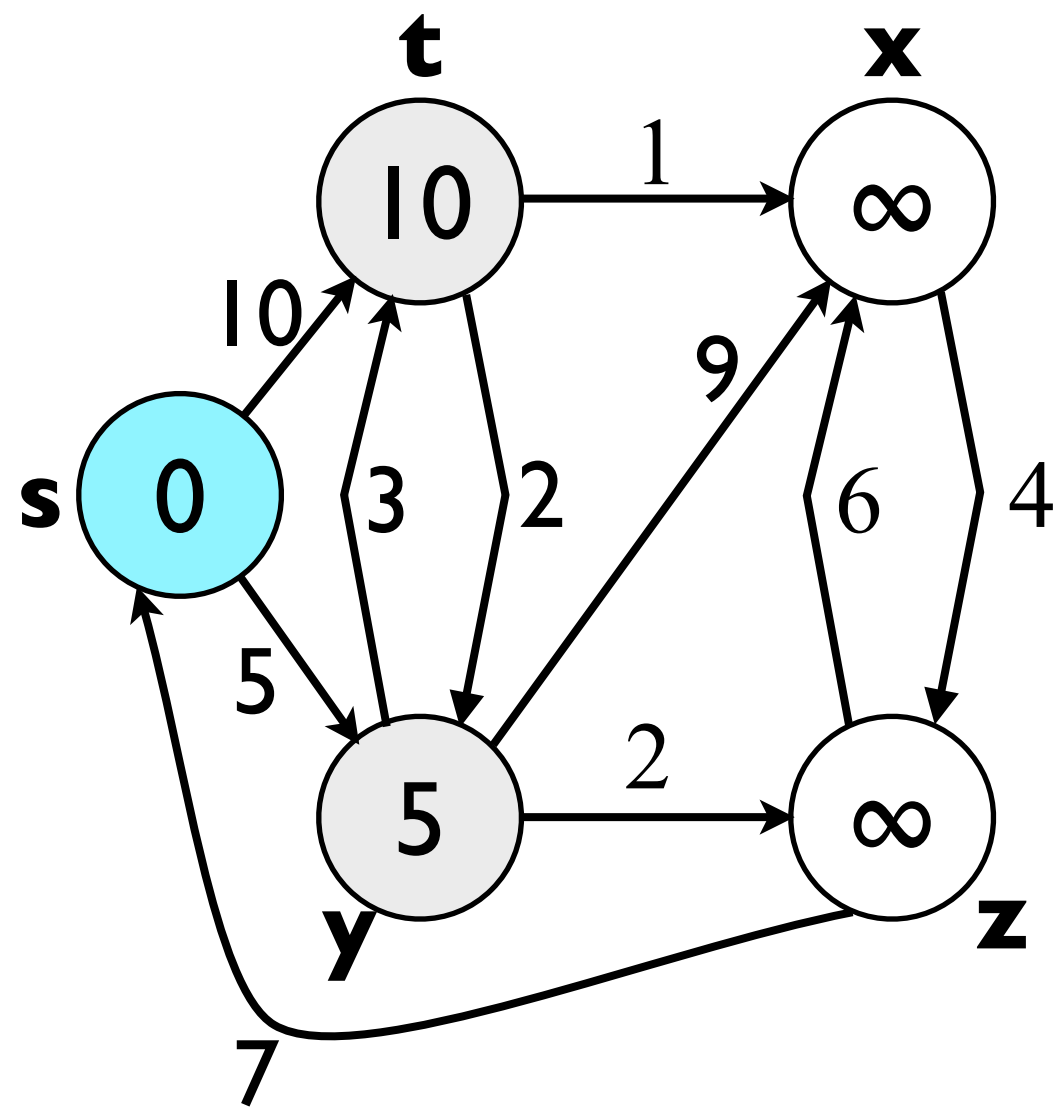
at each step i , finds shortest paths from v_s to nodes of length i

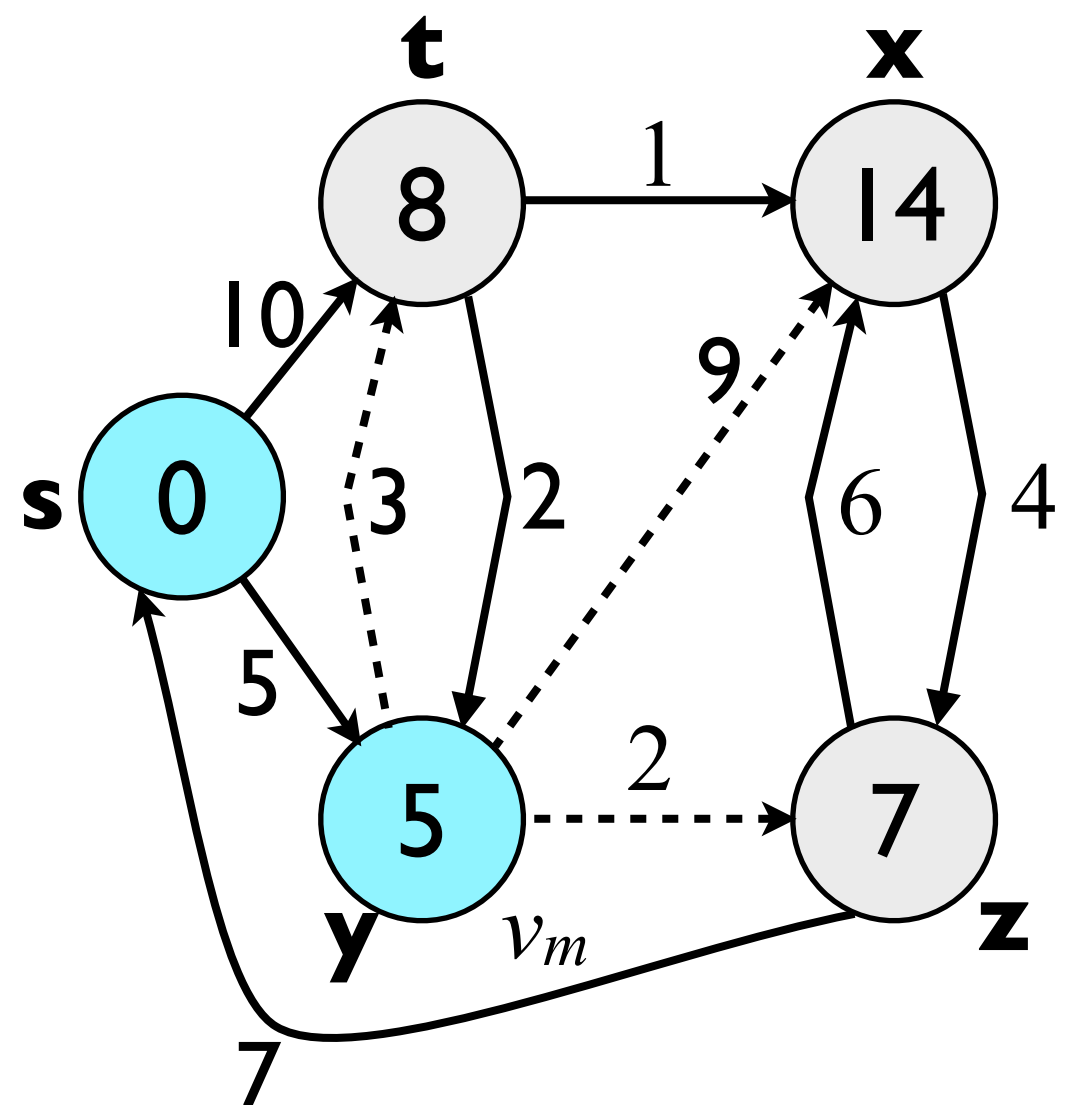
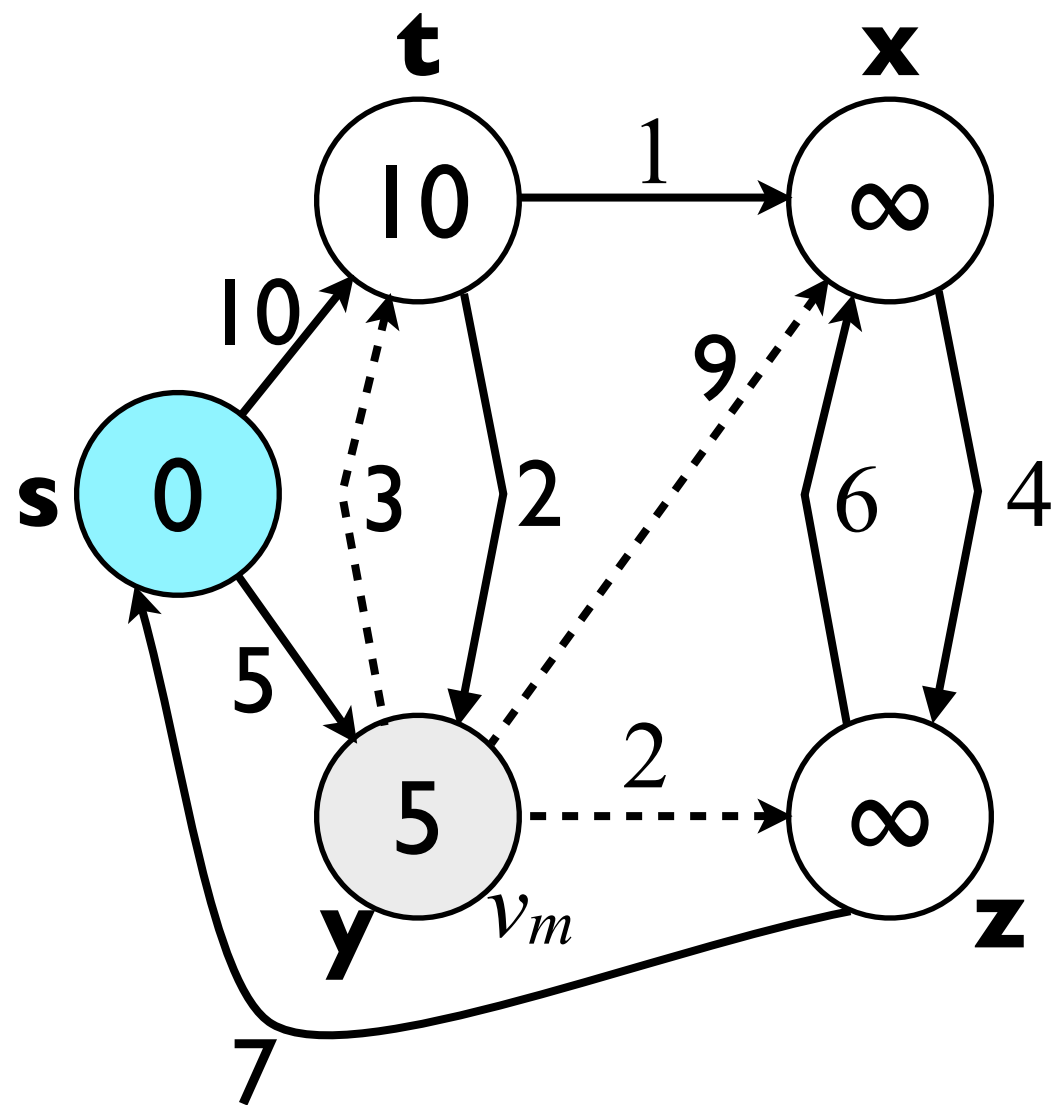


step 1

$$v_m = \mathbf{s}$$

$$d_m = 0$$

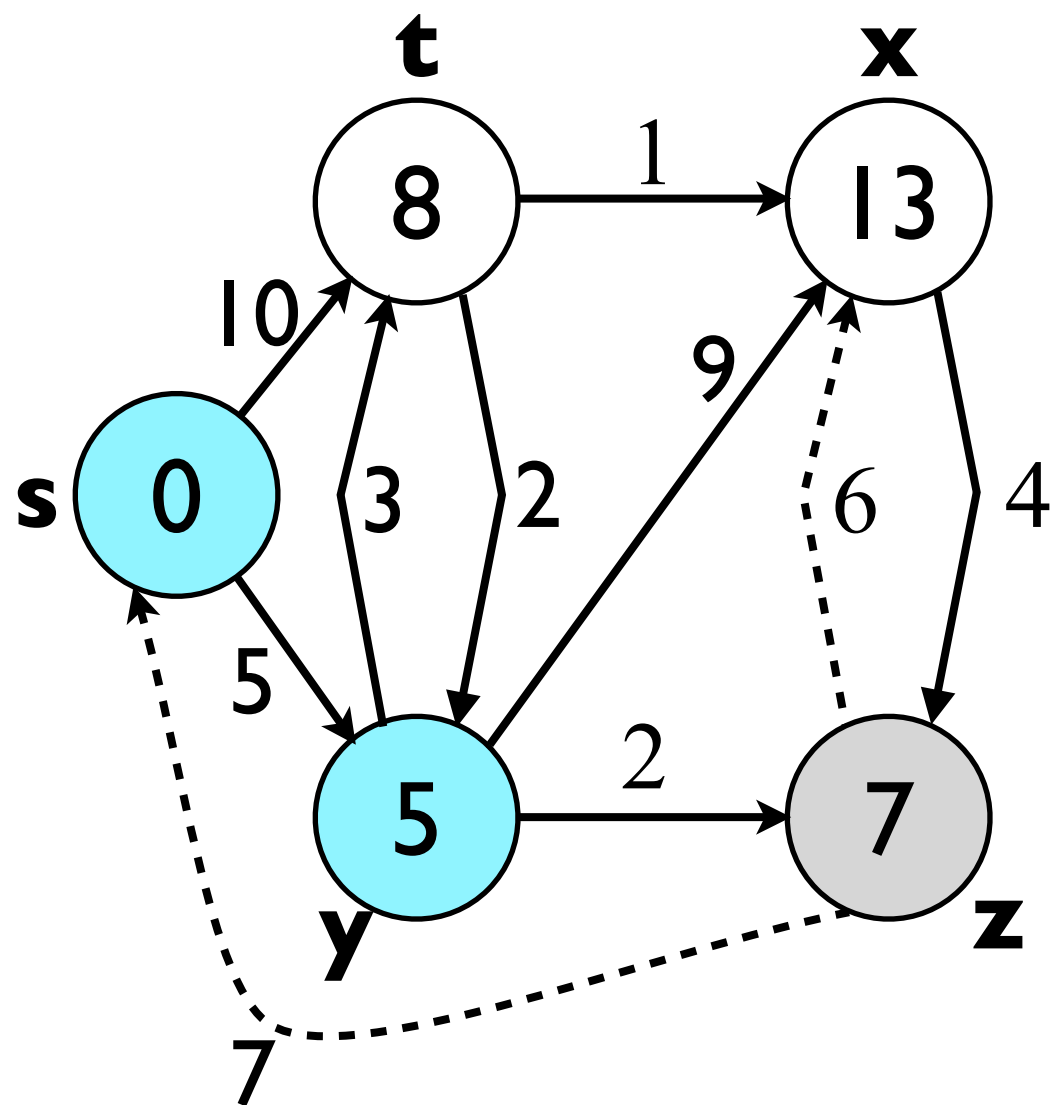




step 2

$$v_m = \mathbf{y}$$

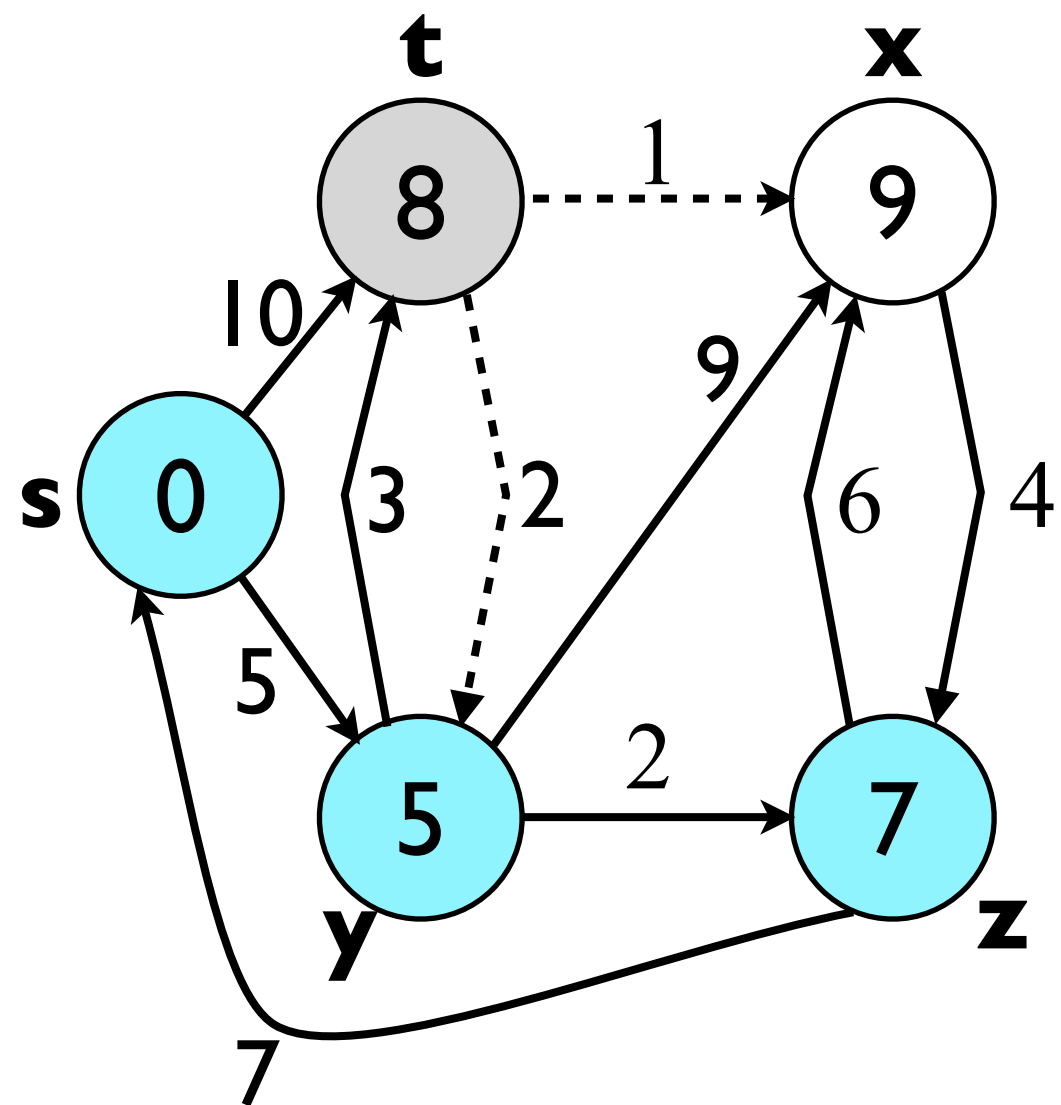
$$d_y = 5$$



step 3

$$v_m = \mathbf{z}$$

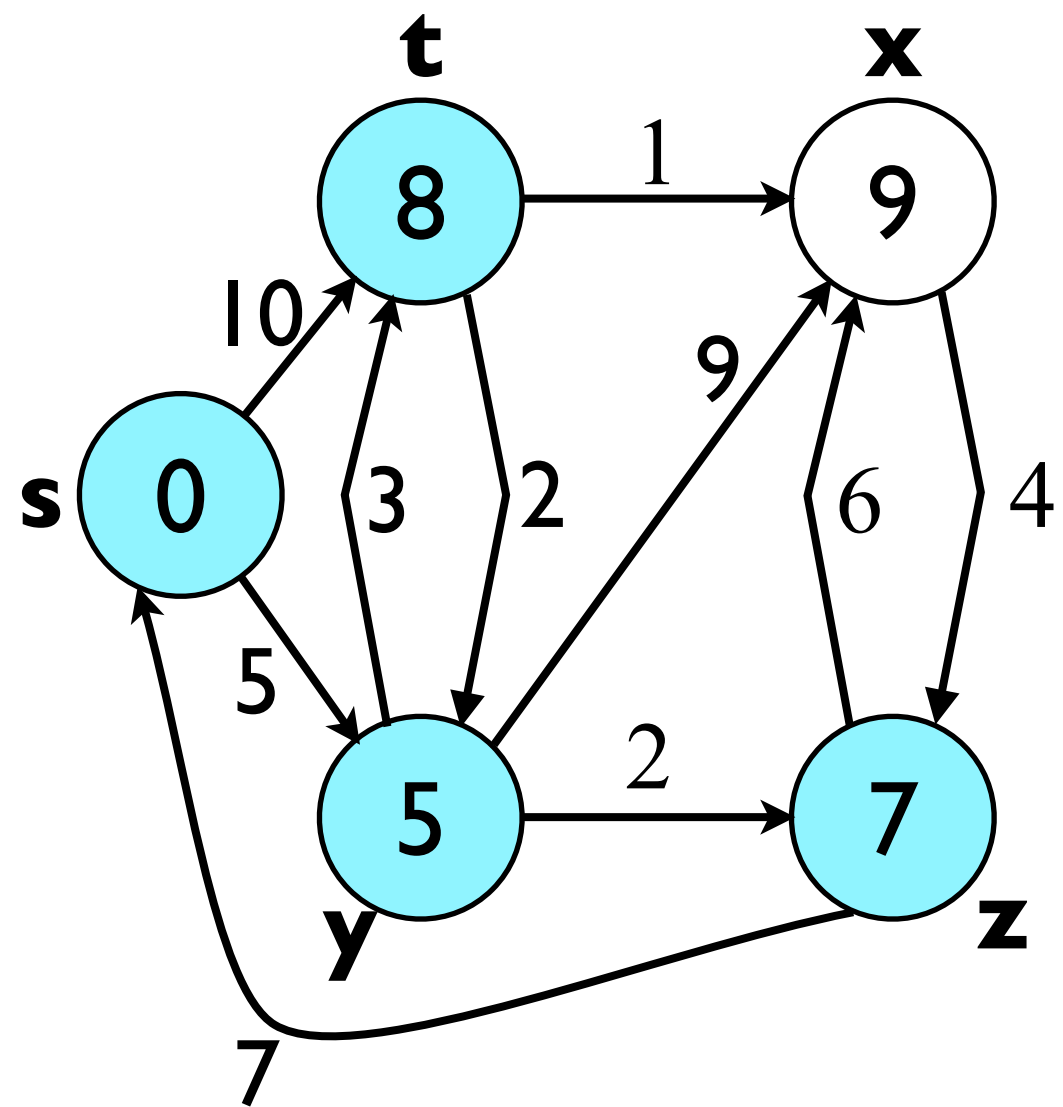
$$d_z = 7$$



step 4

$$v_m = \mathbf{t}$$

$$d_t = 8$$



step 4

$$v_m = \mathbf{t}$$

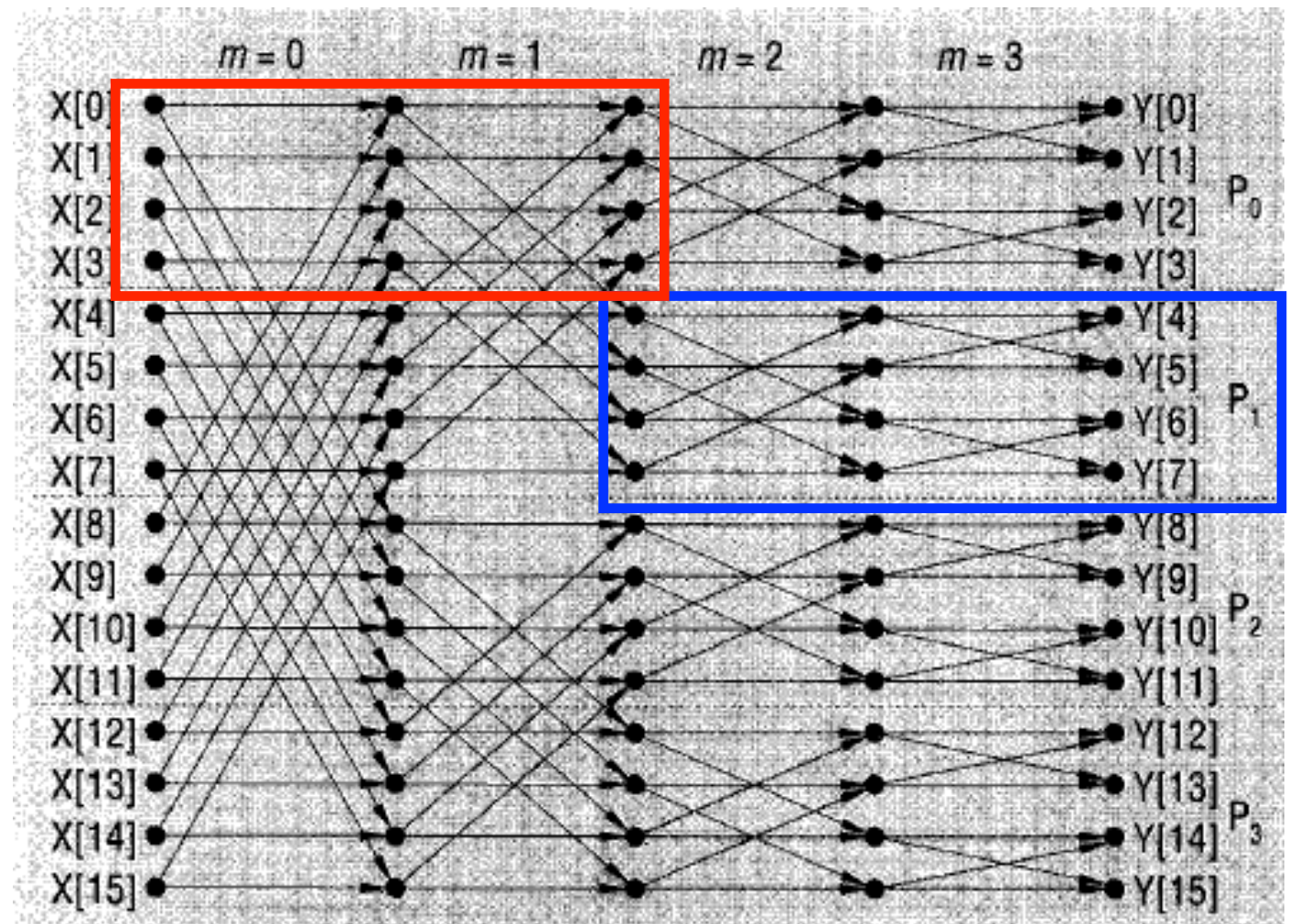
$$d_m = 8$$

A parallel Dijkstra's algorithm for all paths

- Replicate the graph N times (N is the number of vertices), with each processor getting N/P vertices to treat as s vertices, i.e., N/P vertices to find shortest paths from it to other vertices
- Each node computes the shortest distances from the N/P vertices it owns to all other N vertices
- No communication needed
- Seems like the perfect algorithm, but it isn't
- $O(N^3)$ work, but only $O(N)$ parallelism
- W is $\theta(N^3)$, $P=N$, W must grow as $\theta(N^3)$ to scale and the isoefficiency is high

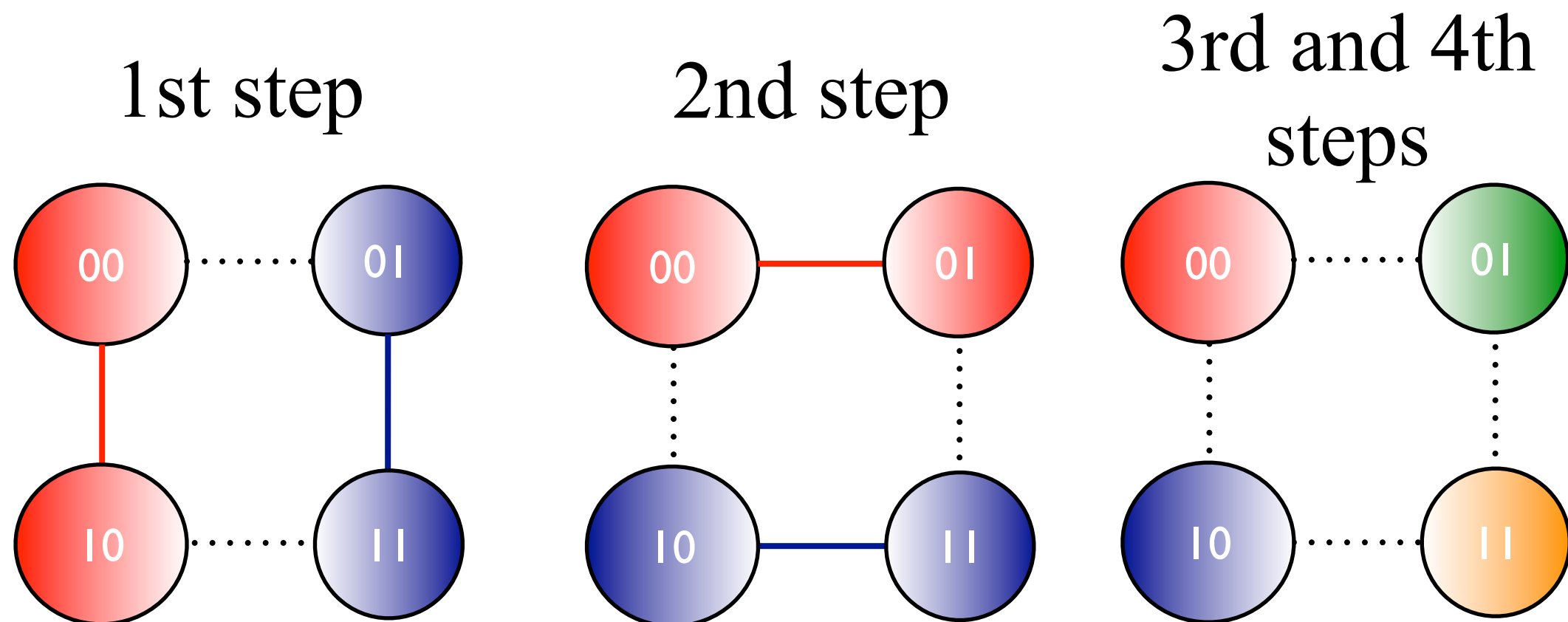
Machine specific parameters

- Consider the Cooley-Tukey FFT algorithm
- Sequential complexity is $\theta(n \log n)$
- Parallel version based on the *binary exchange* method for a d -dimensional ($P=2^d$) hypercube
- partition vectors into block of n/p contiguous elements, $n=2^r$
- 1 block of 2^{r-d} elements assigned per processor
- vector elements on different processors combined during first d iterations, pairs on the same processors combined in the last $r-d$ iterations
- interprocessor communication in only $d = \log P$ of the $r = \log n$ iterations



- Each communication exchanges n/P words
- Communication time is $(t_s + t_w n/P) \log P$
- During each iteration a processor updates n/P elements of vector r
- Let each complex multiply take time t_c

- On a hypercube communicating nodes are *always* adjacent, i.e. a single hop to communicate
- Allows each communication to happen in time $t_s + t_w n/P$ time
- With d communicating steps, hypercube will communicate over each adjacent edge during computation



Parallel execution time

- $T_P = t_c(n/P) \log n + t_s \log P + t_w (n/P) \log P$

computation
time

startup times for $\log p$
communications

$\log P$
communications of n/p
words

- $T_O = P(t_s + t_w n/P) \log P = t_s P \log P + t_w n \log P$
- $W = n \log n$

Solve for different terms

- First term (t_s), $W \propto P t_s \log P$, isoefficiency function is $P \log P$
- Second term, $n \log n = K t_w n \log P$

$$\log n = K t_w \log P$$

$$n = P^{K(t_w/t_c)}$$

$$n \log n = K t_w P^{K t_w} \log P$$

Substituting for K

$$W = E/(1-E) (t_w/t_c) P^{E/(1-E)(t_w/t_c)} \log P$$

Isoefficiency a function of E

- $W = E/(1-E) (t_w/t_c) P^{E/(1-E)(t_w/t_c)} \log P$
- Consider if exponent of P, $t_w E/(t_c (1-E)) < 1$

W grows slower than $P \log P$

Overall isoefficiency is $\theta(P \log P)$ (from t_s term prev. page)

- Consider if $t_w E/(t_c (1-E)) > 1$

isoefficiency a function of relative values of $E/(1-E)$, t_w , t_c

- Consider if $t_w E/(t_c (1-E)) = 1$

Isoefficiency is $P \log P$, a lower threshold for a hypercube

effect of $t_w E/(t_c (1-E))$ on isoefficiency

- if $t_w = t_c$, isoefficiency is $W = E/(1-E) P^{E/(1-E)} \log P$
- Now for $E/(1-E) \leq 1, E \leq 0.5$
isoefficiency is $\theta(P \log P)$
- For $E/(1-E), E \geq 0.5$
If $E=0.9, E/(1-E)=9$, isoefficiency is $P^9 \log P$
- Effect of t_w and t_c : let's make the bandwidth lower
if $t_w = 2t_c$ then the *threshold* efficiency is 0.333
Isoefficiency for $E=0.333$ is $\theta(P \log P)$
Isoefficiency for $E=0.5$ is $\theta(P^2 \log P)$ and for
 $E=0.9$ is $\theta(p^{18} \log p)$ ($t_w E/(t_c (1-E)) = 2E(1-E) = 1.8/0.1$)
- What can we conclude from this?

Conclusions for FFT

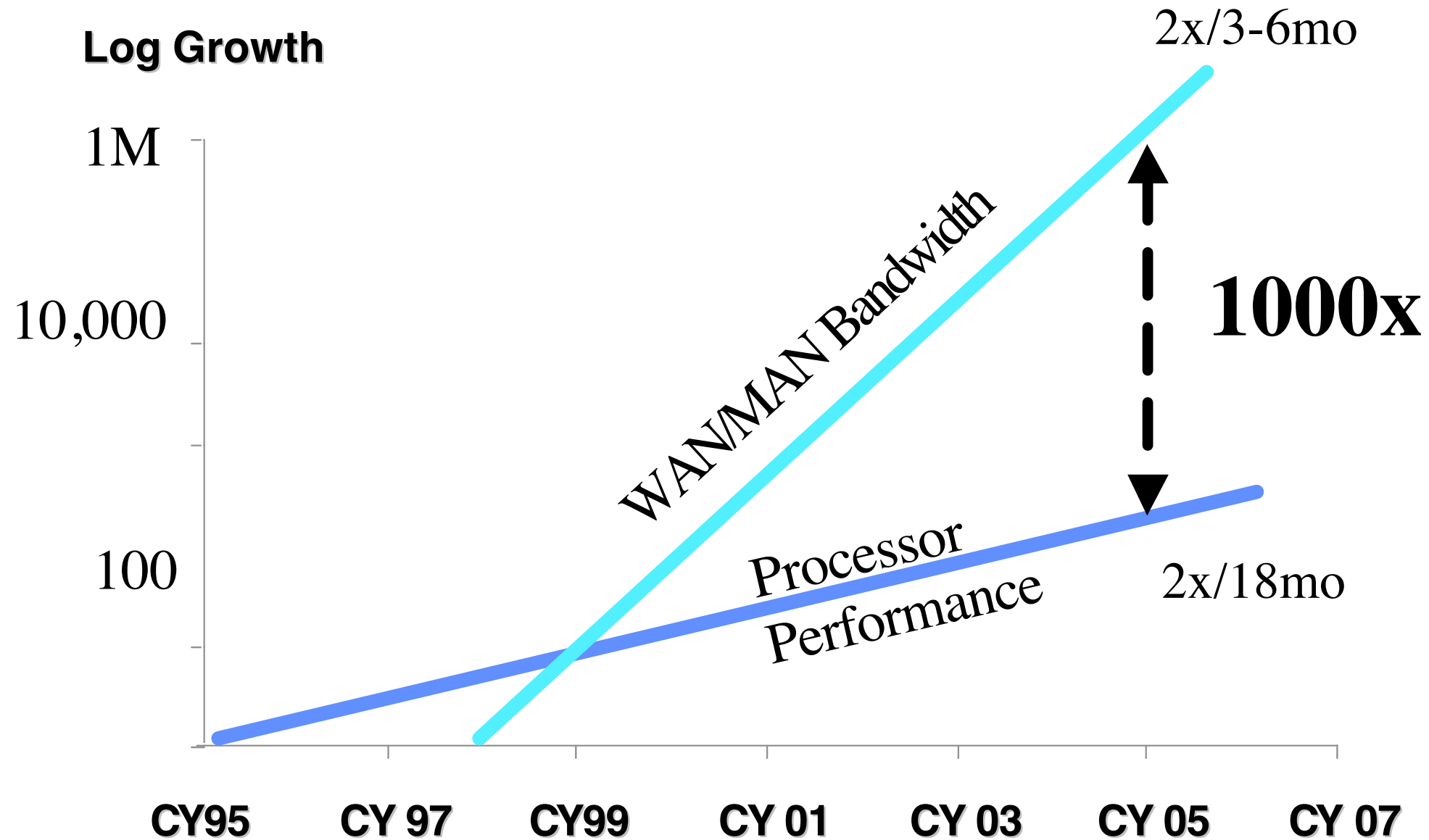
- Balance of bandwidth and CPU is important for this problem - scalability is good on a balanced system
- Making bandwidth higher helps
- Increasing CPU performance without increasing bandwidth reduces scalability
- On modern systems ...

From a talk by Horst Simon

Bandwidth vs. Moore's Law



Adapted from G. Papadopoulos, Sun



FFT is unique in this property

- **But**, the ratios of t_w and t_c can be high
- May result in t_c term being important in small machine sizes, and the t_w or t_s terms dominating for larger machines
- Again, need to apply intelligence, and again, using isoefficiency gives insights into what is required to have an app scale

Summary

- Data structure contention also must be considered if it is the dominating term
- In summary:
 - Want to increase problem size to maintain efficiency
 - Must have enough memory to hold larger problem size
 - Rate of growth of problem size is a limit on the number of processors we can run on
 - Thus rate of growth of problem size is a limit on how scalable the algorithm is *if we want to maintain constant efficiency*
 - Isoefficiency functions provide a way of determining the rate of growth of the problem size