# High Performance Switching and Routing
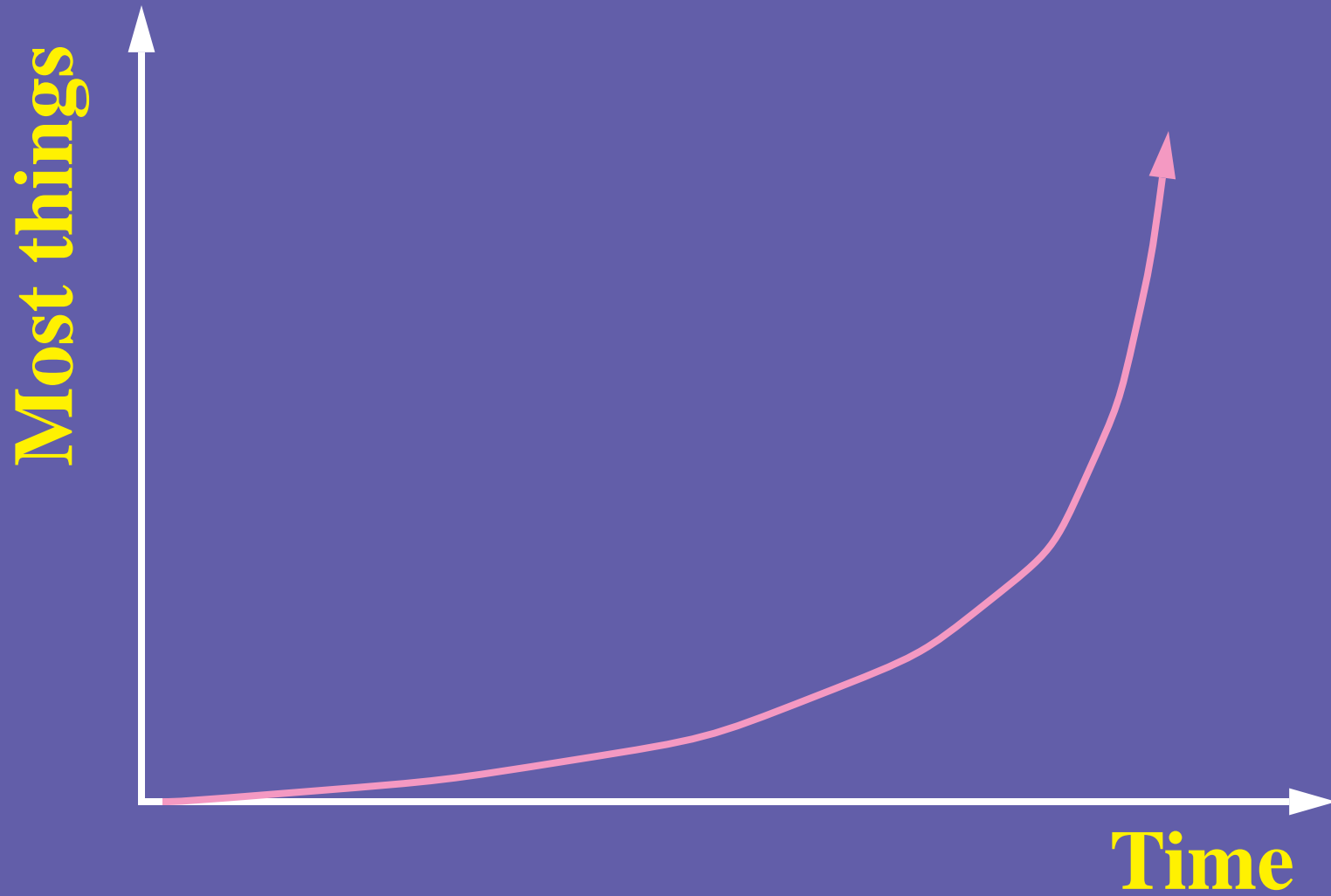
**Nick McKeown**

**Assistant Professor of Electrical Engineering and Computer Science**

nickm@stanford.edu
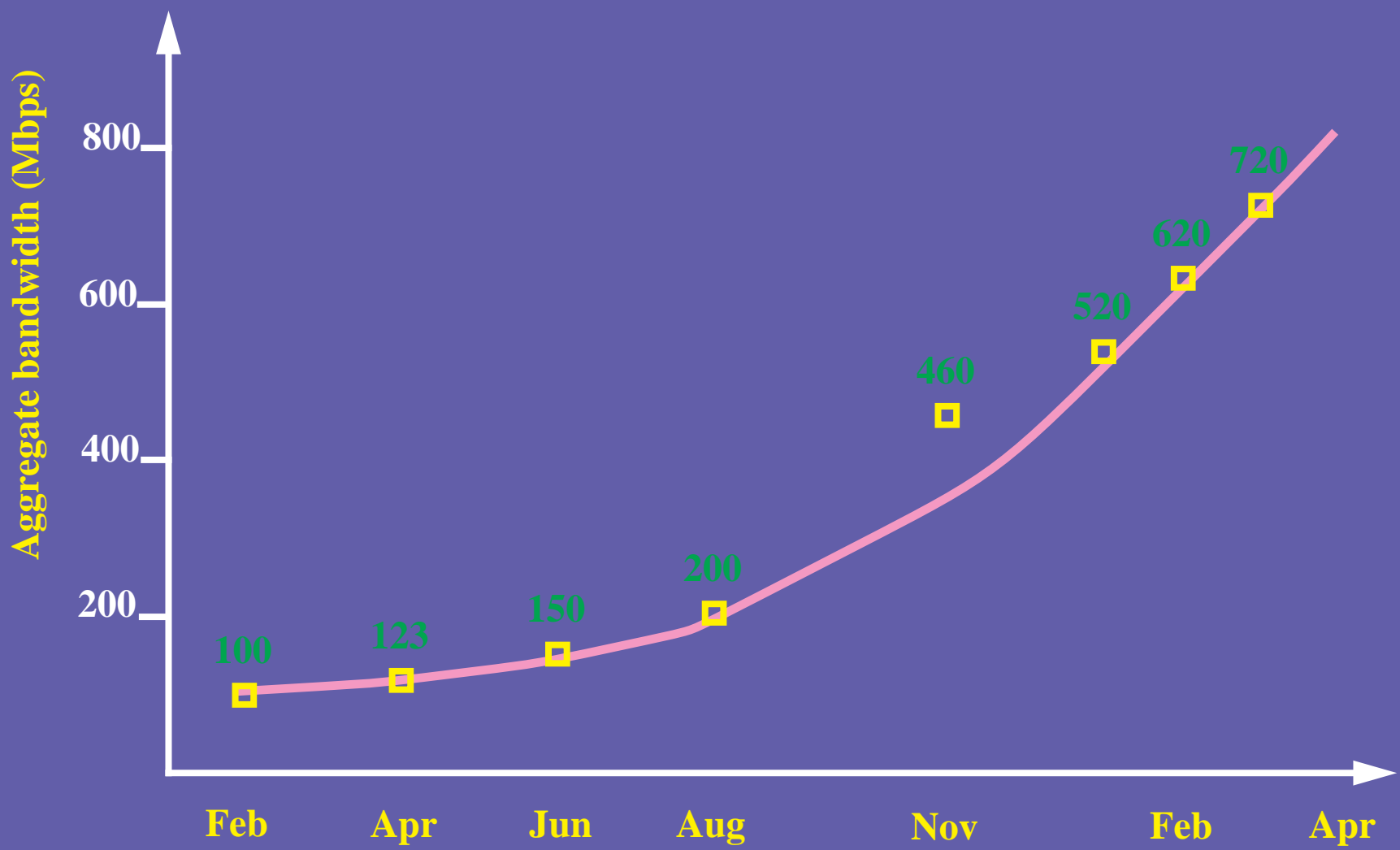http://www.stanford.edu/~nickm

1. The Demand for Bandwidth

2. The Shortage of Switching/Routing Capacity

3. The Architecture of Switches and Routers
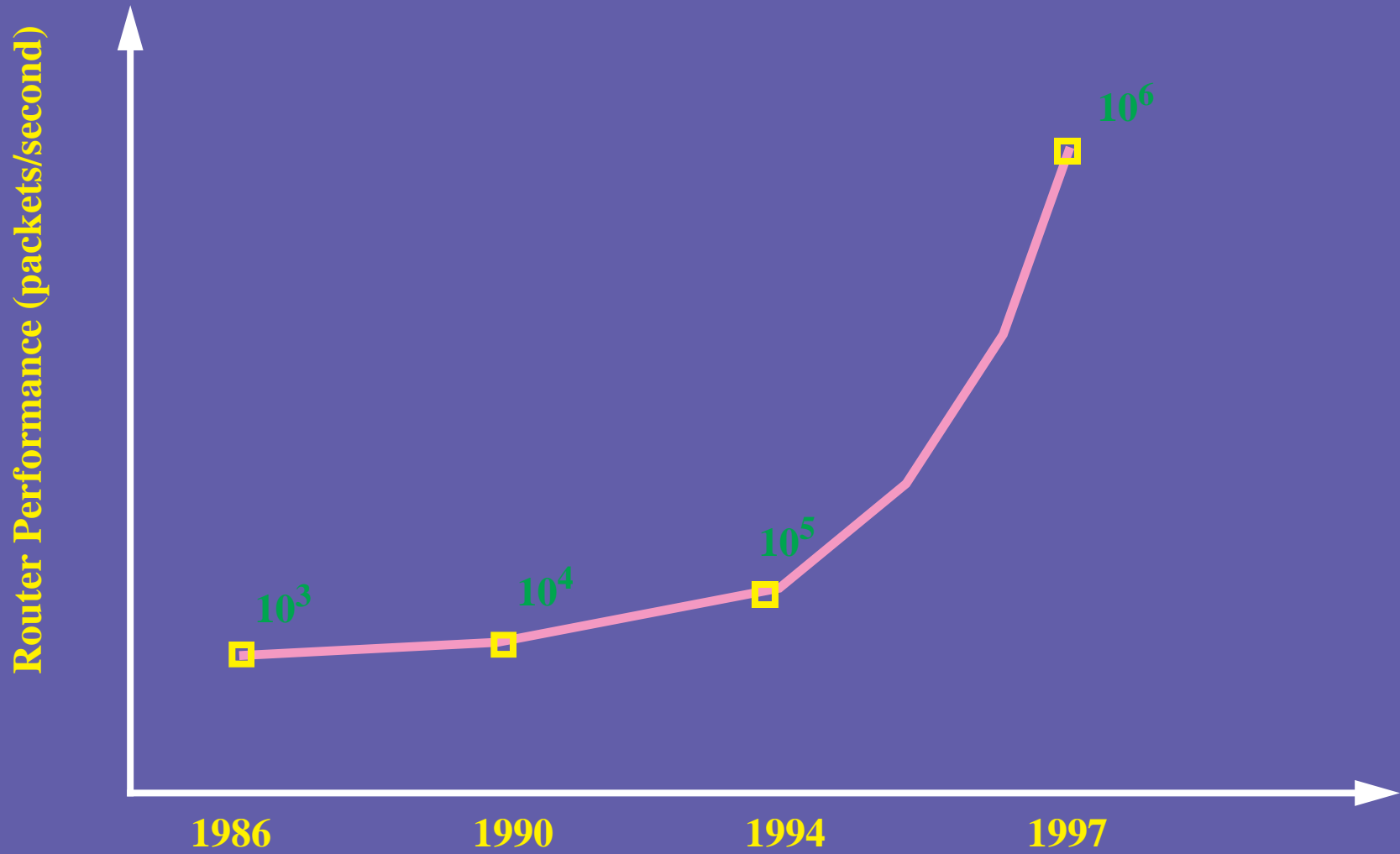
4. Some (of our) solutions

# What's the Problem?



Most things vs. Time
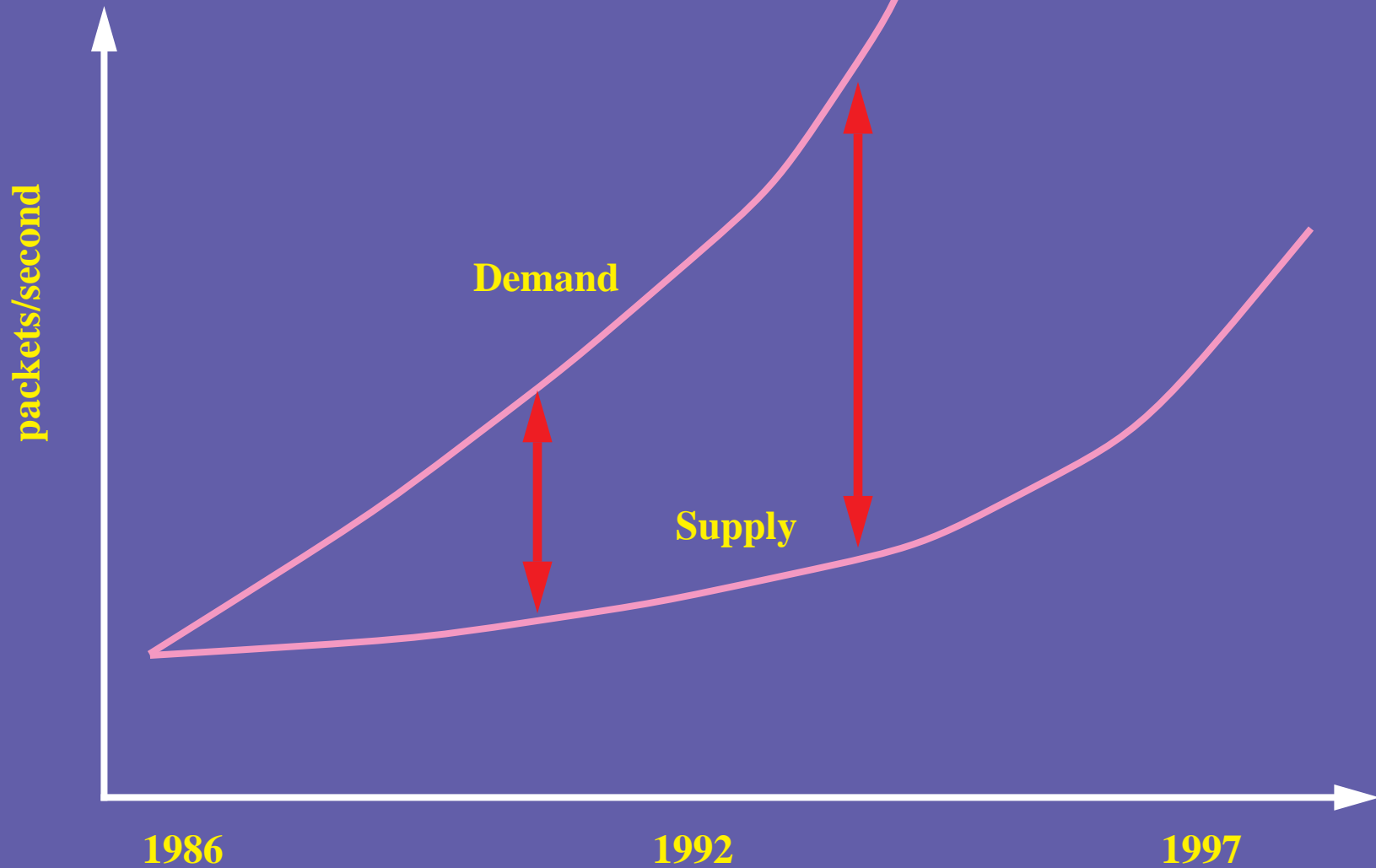
# The demand
## *The San Jose NAP*



Aggregate bandwidth (Mbps)

100  123  150  200  460  520  620  720

Feb  Apr  Jun  Aug  Nov  Feb  Apr

*Source: http://www.mfsdatanet.com/MAE/west.stats.html*

High Performance Switching and Routing

4

# The supply



Router Performance (packets/second)

$10^6$

$10^5$

$10^4$

$10^3$

1986    1990    1994    1997

# Why we need faster switches/routers



packets/second

Demand

Supply

1986            1992            1997

High Performance Switching and Routing

6

# Traffic Inversion
## *10 years ago*

## *Today*

# Why is this a problem?



11/01/96 Packet Loss for BBNPlanet (AS1) Between Mae-West and Sprint

Packet Loss (%)

mae-west:sprint

**November 1st, 1996**

1. The Demand for Bandwidth

2. The Shortage of Switching/Routing Capacity

☞ 3. The Architecture of Switches and Routers

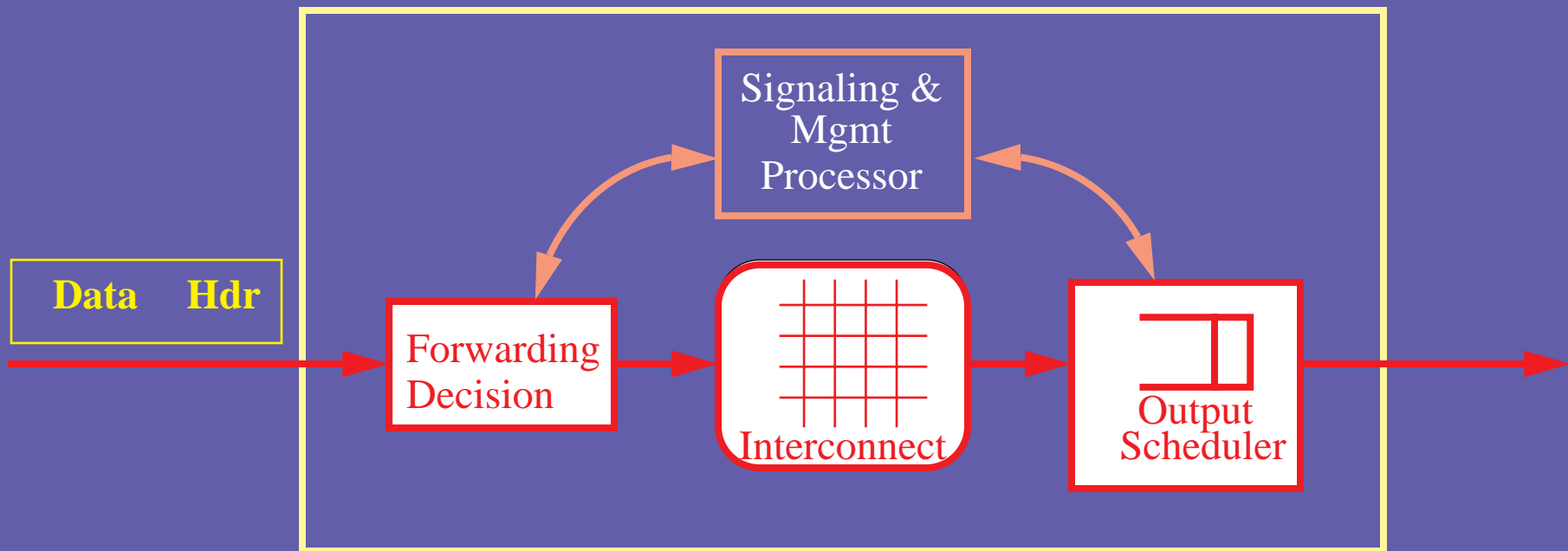4. Some (of our) solutions

# The Architecture of Switches and Routers

Generic Packet Switch:
(e.g. IP Router, ATM Switch, LAN Switch)



Signaling & Mgmt Processor

Data    Hdr

Forwarding Decision

Interconnect

Output Scheduler

# Performance of IP Routers



Time

Copy
Time

Forwarding
Decision
Time

Arrival
Time

Min back-to-back packet size

Packet size

# Performance of IP Routers



**Most routers do this poorly!**

**Time**

Copy rate

Header processing time

Arrival rate

Min back-to-back packet size

**Packet size**

**Most routers do this ~ ok**

# The Evolution of Routers

## *The first shared memory routers*

**Routing CPU**

**Buffer Memory**

**DMA**
**Line Card**
**MAC**
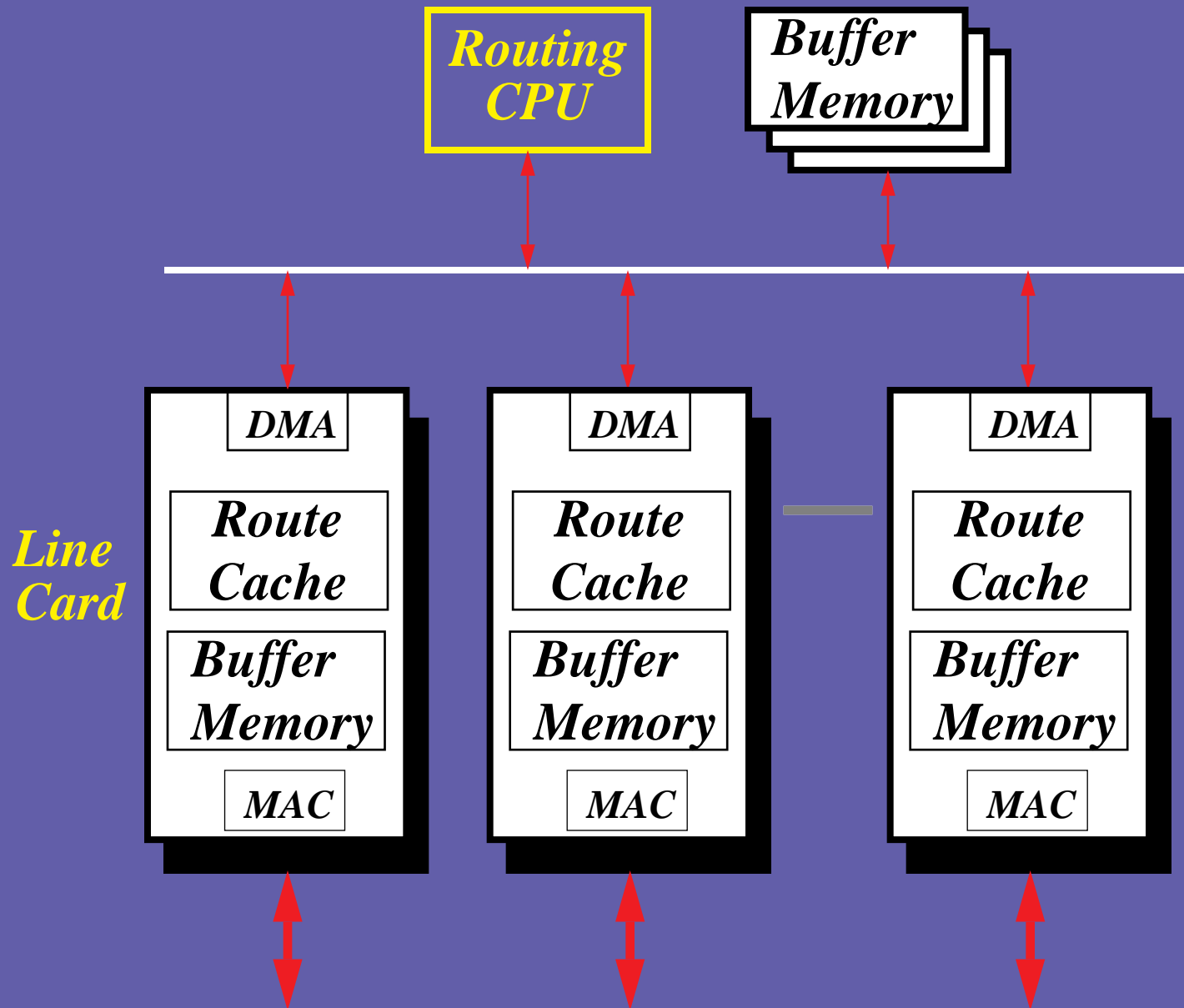
**DMA**
**Line Card**
**MAC**

**DMA**
**Line Card**
**MAC**

# The Evolution of Routers
## *The first shared memory routers*

# The Evolution of Routers

## *Reducing the number of bus copies*



Routing CPU

Buffer Memory

Line Card

DMA — Route Cache — Buffer Memory — MAC

DMA — Route Cache — Buffer Memory — MAC

DMA — Route Cache — Buffer Memory — MAC

# The Evolution of Routers

*Reducing the number of bus copies*

**updates**

**Routing CPU**

**Buffer Memory**

**Line Card**

| DMA | DMA | DMA |
| Route Cache | Route Cache | Route Cache |
| Buffer Memory | Buffer Memory | Buffer Memory |
| MAC | MAC | MAC |

# The Evolution of Routers

## *Avoiding bus contention*

**Memory Buffer**

**ROUTE CPU**

**Route Cache**

**Buffer Memory**

**MAC**

**Route Cache**

**Buffer Memory**

**MAC**

**Route Cache**

**Buffer Memory**

**MAC**

**Advantage:**
Non-blocking backplane— high throughput
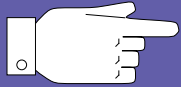
**Disadvantage:**
Difficult to provide QoS

1. The Demand for Bandwidth

2. The Shortage of Switching/Routing Capacity

3. The Architecture of Switches and Routers

☞ 4. Some (of our) solutions
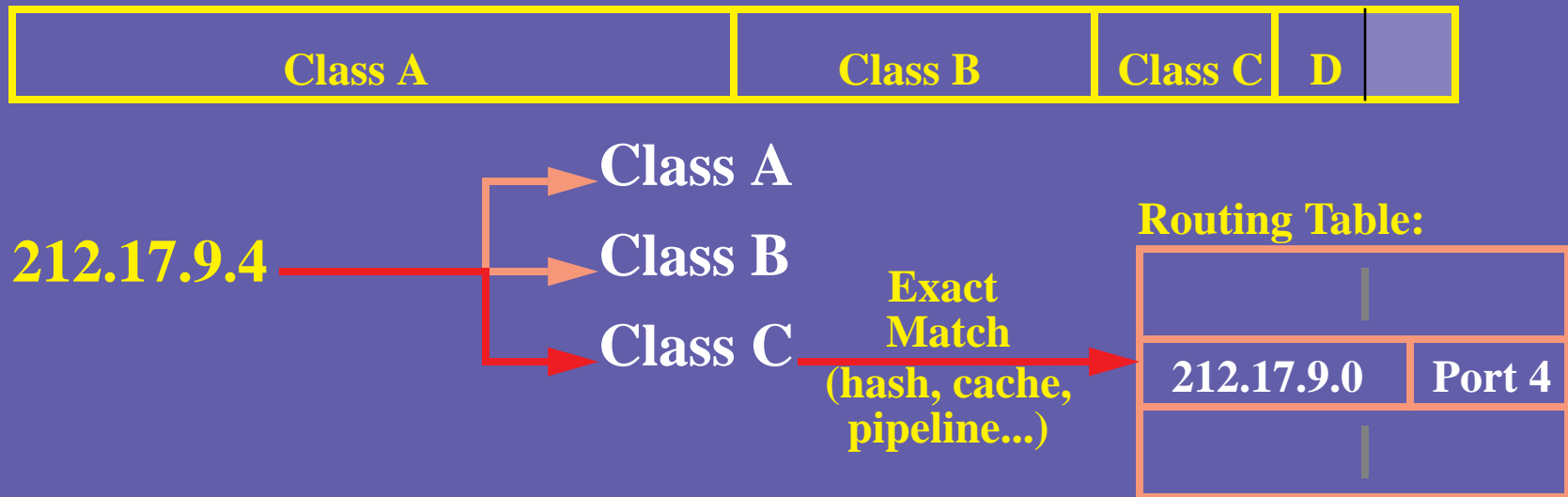
# Some (of our) Solutions

**1. Accelerating Forwardng Decisions:**

- Longest-matching prefixes
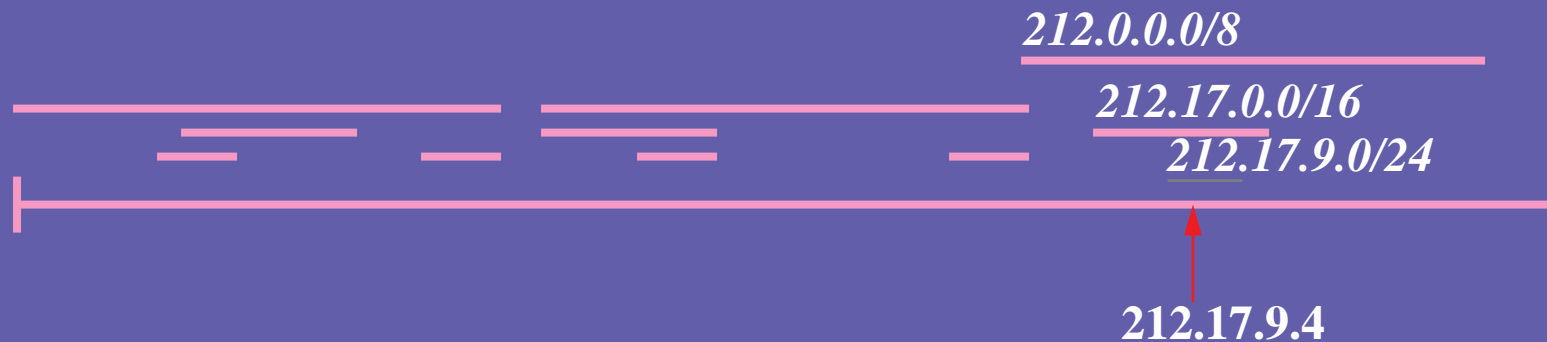
**2. Interconnections: Switched Backplanes**

- Input Queueing
  — Theory
  — Unicast
  — Multicast

- Fast Buffering

- Speedup

- The *Tiny Tera* Project

# Routing Lookups

| Class A | Class B | Class C | D | |

212.17.9.4 → Class A / Class B / Class C

**Class A**

**Class B**

**Class C** → **Exact Match (hash, cache, pipeline...)** →

**Routing Table:**

| | |
|---|---|
| | |
| 212.17.9.0 | Port 4 |
| | |

# Routing Lookups with CIDR ("supernetting")

**CIDR uses "longest matching prefix" routing:**

*212.0.0.0/8*

*212.17.0.0/16*

*212.17.9.0/24*

**212.17.9.4**

*Hashing, caching and pipelining are hard!*

# Perform Lookups Faster

*Observation #1:*



*Size of Routing Tables*

*Cost of Memory (per byte)*

*Time*

# *Performing Lookups Faster*

*Observation #2:*

*Number in routing table*

*Prefix length*

*24*

*256*

*212.17.9.0/24*

*0*

*212.17.9.4*

$2^{32}$-1

# 20 million lookups per second

**16Mbytes of 50ns DRAM**

**212.17.9.1**

| | |
|---|---|
| **1** | **Port 4** |
| **0** | **look further** |
| **1** | **Port 4** |
| **1** | **Port 3** |
| | |
| **0** | **look further** |
| **1** | **Port 3** |

**<1Mbyte of 50ns DRAM**

| |
|---|
| *Port 4* |
| *Port 5* |
| |

**256**

| |
|---|
| *Port 4* |
| *Port 5* |
| |

**1.  Accelerating Forwardng Decisions:**

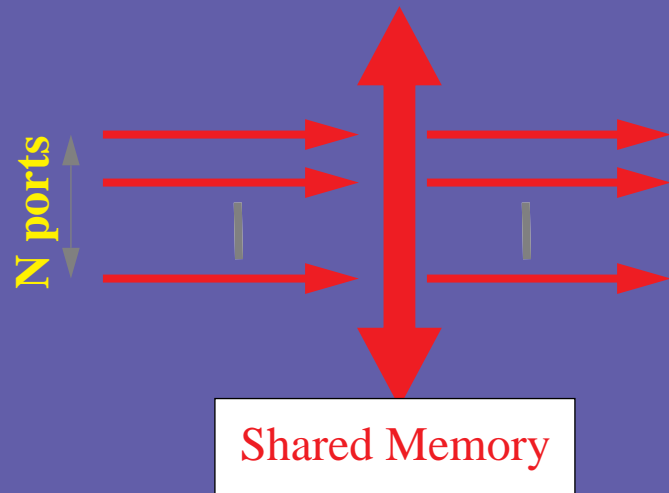- Longest-matching prefixes

**2. Interconnections: Switched Backplanes**

- Input Queueing
  - — Theory
  - — Unicast
  - — Multicast

- Fast Buffering

- Speedup

- The *Tiny Tera* Project

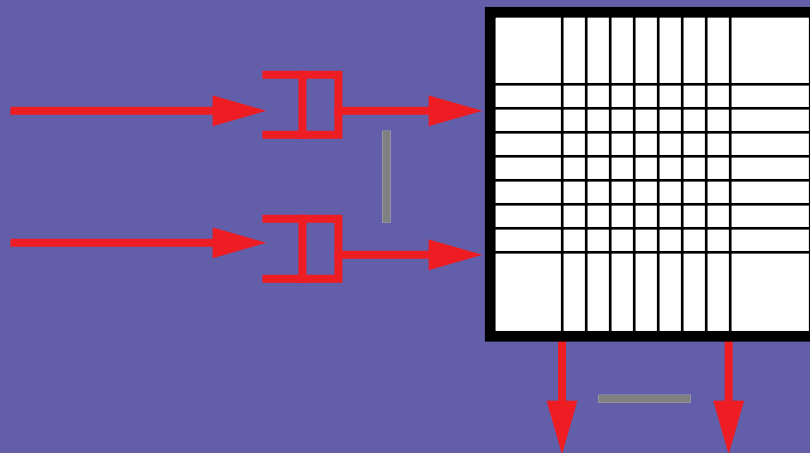# Should we use shared memory or input-queueing?

## Shared Memory:



N ports

Shared Memory

**Advantages:**

Highest Throughput.
Possible to control packet delay.

**Disadvantages:**

N-fold internal speedup
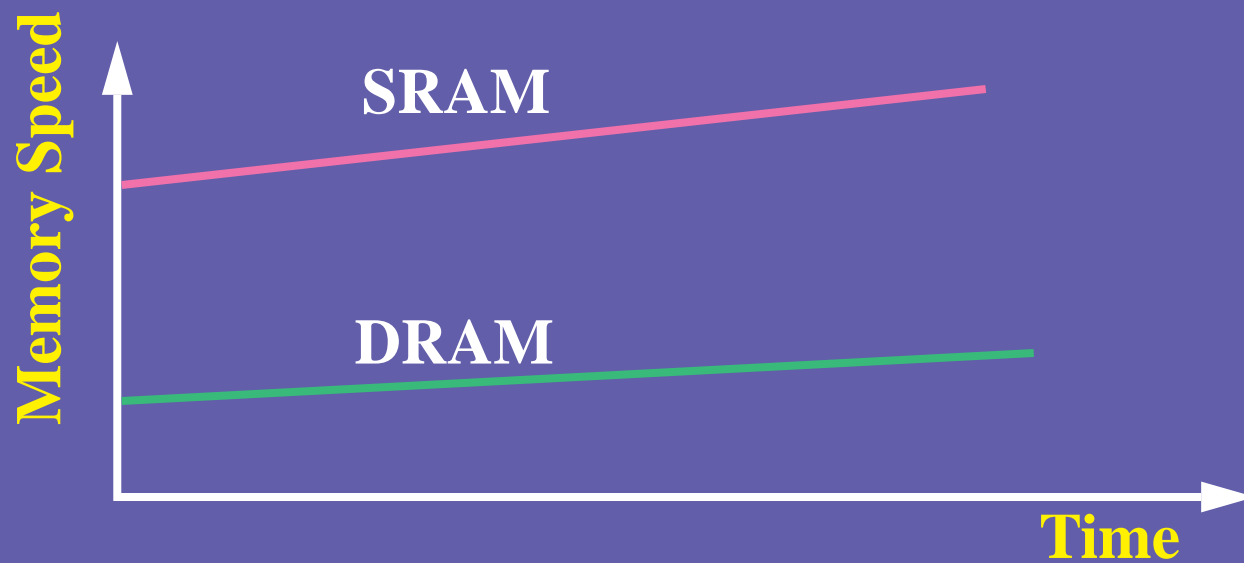
## Input Queueing:
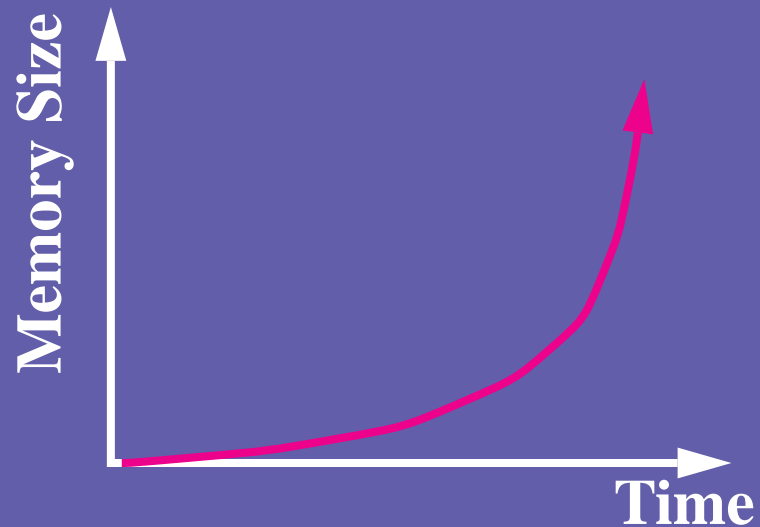


**Advantages:**

Simplicity
High Bandwidth

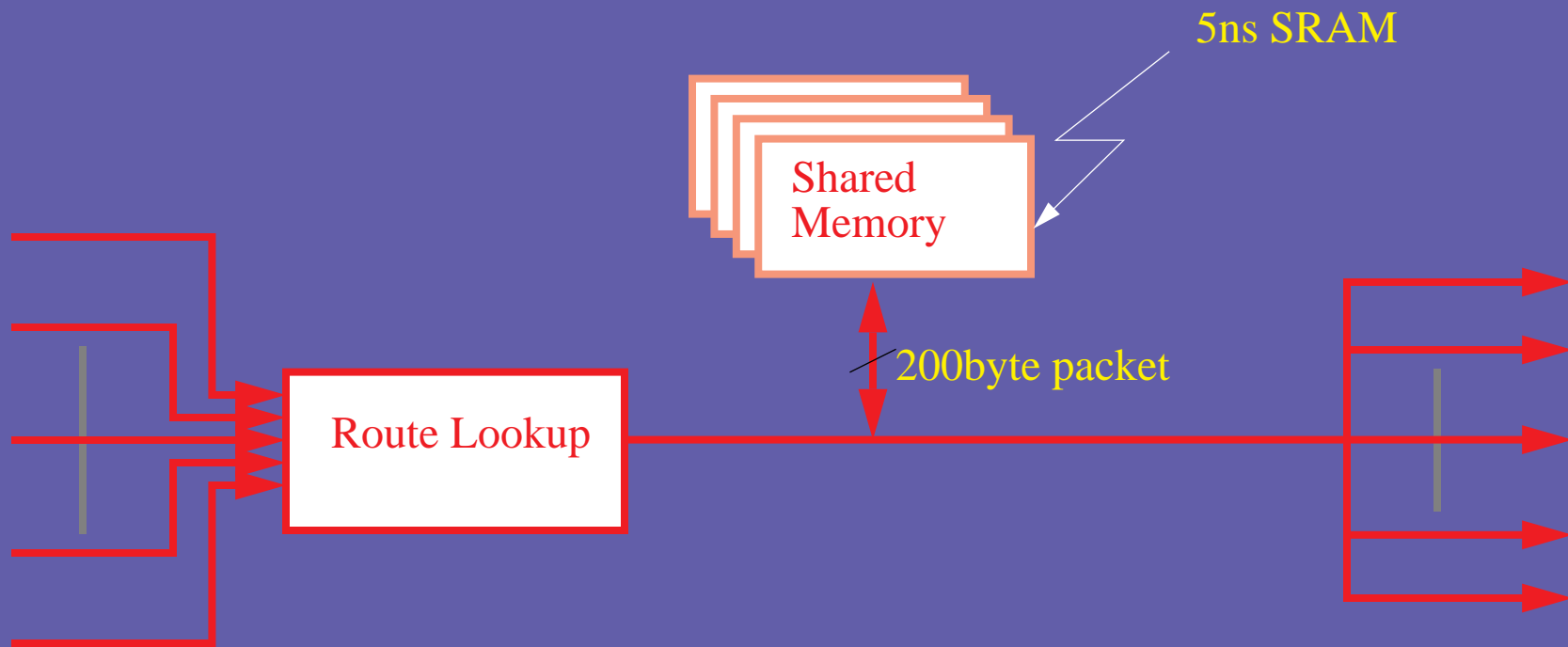**Disadvantages:**

HOL Blocking
Less efficient
Difficult to control packet delay.

# Memory Bandwidth

# *An aside:* How fast can shared memory operate?

5ns SRAM

Shared Memory

200byte packet

Route Lookup

## *How fast can a 16 port switch run with this architecture?*

*5ns per packet × 2 memory operations per cell time*

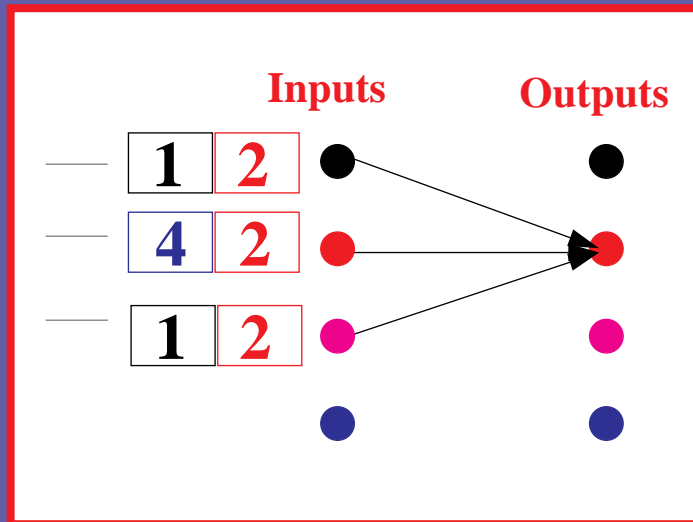$\Rightarrow$ *aggregate bandwidth is 160Gb/s*

# Should we use shared memory or input-queueing?

Because of a *shortage of memory bandwidth*, most multigigabit and terabit switches and routers use either:

1. Input Queueing, or
2. Combined Input and Output Queueing.
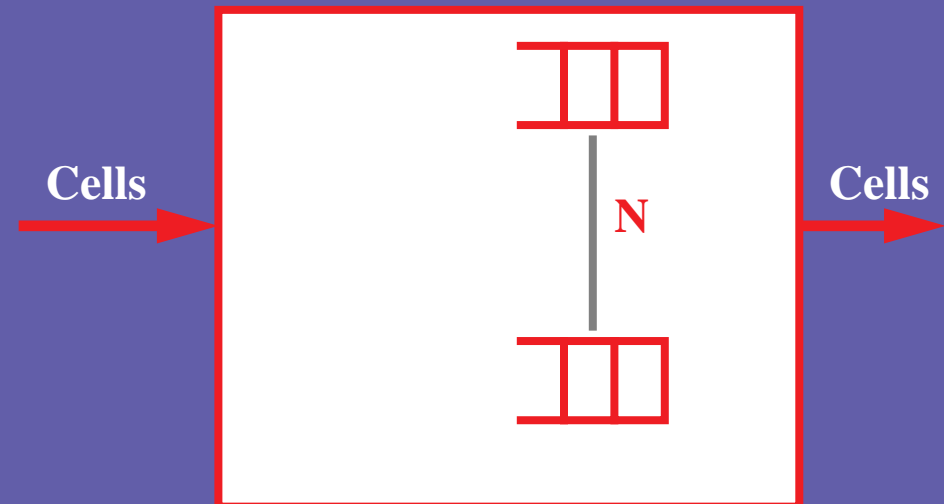
# Head of Line Blocking

## The Problem



$$\rho_{max} = 2 - \sqrt{2} = 58\%$$

## A Solution....

### Input Cell Buffer



*"Virtual Output Queueing"*

$$\rho_{max} = 100\%$$

# ...but requires scheduling...

Input 1

$Q(1,1)$

$A_{1,1}(t)$

$A_1(t)$

$Q(1,n)$

Matching, M

Output 1
$D_1(t)$

?

Input m

$Q(m,1)$

$A_m(t)$

$Q(m,n)$

Output n
$D_n(t)$

# ....which is equivalent to graph matching

**Request Graph**

1 — 7 — 1
2
4
2
3
5
4 — 2

**Bipartite Matching**

(Weight = 18)

# Practical Algorithms

1. *i*SLIP  — **Weight = 1**

          — **Iterative round-robin**

          — **Simple to implement**

*Simple, fast, efficient*

2. *i*LQF  — **Weight = Occupancy**

3. *i*OCF  — **Weight = Cell Age**

*Good for non-uniform traffic. Complex!*

4. LPF   — **Weight = Backlog**

*Good for non-uniform traffic. Simple.*

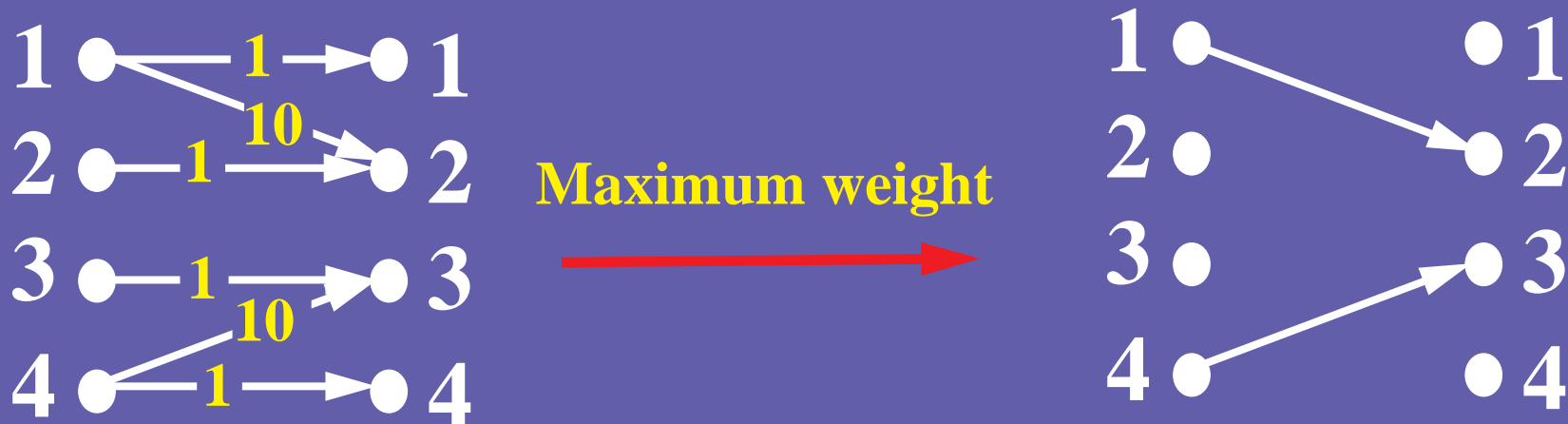# Achieving 100% Throughput
## Longest Queue First & Oldest Cell First

$$\text{Weight} = \left\{ \begin{array}{l} \text{Queue Length} \\ \text{Waiting Time} \end{array} \right\} \Rightarrow 100\%$$

Maximum weight

# Theorems

## Theorem:

*Both LQF and OCF can achieve 100% throughput for independent traffic both uniform and non-uniform.*

## Proof:

$$E\left[\sum_{i,j} L_{i,j}(n)\right] < \infty, \ \forall n \quad \overset{\textbf{Def}}{\Longrightarrow} \quad \textbf{100\% throughput}$$
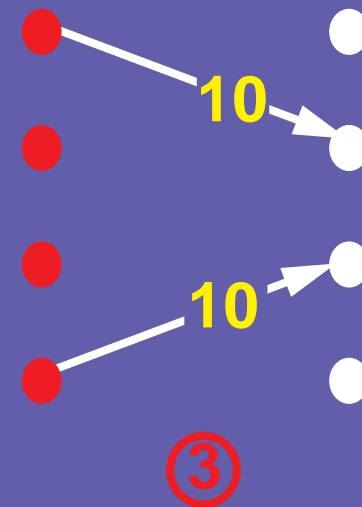
$\Uparrow$

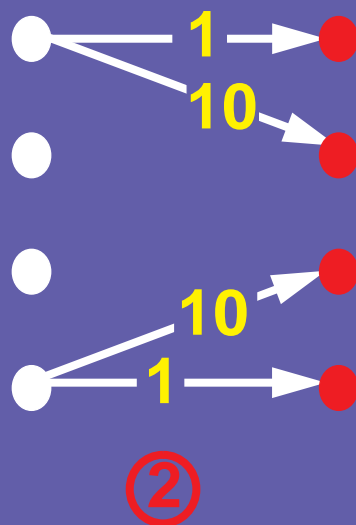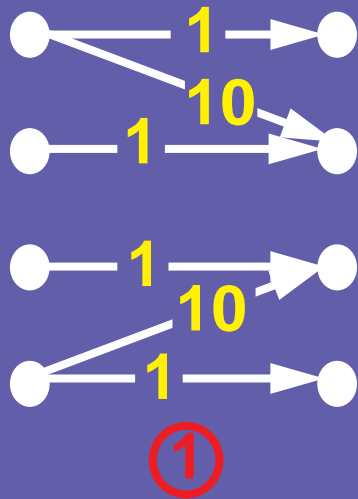**Lyapunov Stability Criterion:**

$$E\left[V(\underline{L}(n+1)) - V(\underline{L}(n)) \mid \underline{L}(n)\right] \leq 0, \forall |\underline{L}(n)| > k$$

http://tiny-tera.stanford.edu/~adisak/research.html

## iLQF&iOCF



**Iteration steps**

Step 1. *Request*

Step 2. *Grant* to the largest request

Step 3. *Accept* grant to the largest request

# *i*LQF and *i*OCF
# Problem is in Comparators

# Solution to
## *Complexity Problem*

☞ **Longest Port First (LPF)**

☞ **Oldest Port First (OPF)**

**Advantages**

— **SIMPLER.**

  • **Can use maximum size matching — O($N^{2.5}$).**

— **FASTER.**

  • **Move magnitude comparator out of the critical path.**

  • **Lends itself well to pipelining.**

# LPF Algorithm
# Using Port Occupancy



$$w_{i,j} = \sum_{j} L_{i,j} + \sum_{i} L_{i,j}$$

i.e. $w_{1,1} = L_{1,1} + L_{1,2} + L_{1,1} + L_{2,1}$

**Input occupancy**     **Output occupancy**

# On The Theorems

## Theorem:

*An LPF match is of both maximum weight and maximum size.*

## Theorem:

*LPF can achieve 100% throughput for independent traffic both uniform and non-uniform.*

**Proof:**  $V(\underline{L}(n)) = \underline{L}^T(n)T\underline{L}(n)$

$$E\left[V(\underline{L}(n+1)) - V(\underline{L}(n)) \mid \underline{L}(n)\right] \leq 0, \forall \left|\underline{L}(n)\right| > k$$
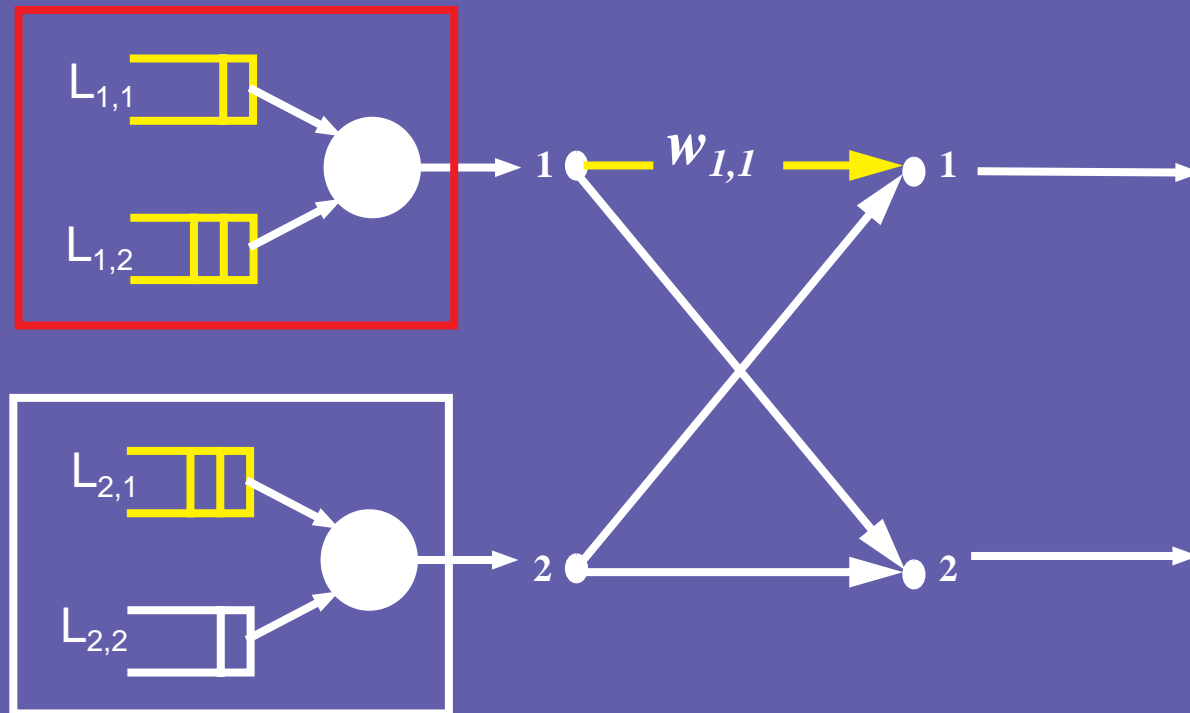
# Presorting Inputs & Outputs

$$\begin{bmatrix} 2 & 20 & 0 \\ 17 & 0 & 8 \\ 0 & 0 & 1 \end{bmatrix}$$

*Weight request* →

|      | **19** | **20** | **9** |
|------|--------|--------|-------|
|      | **1**  | **1**  | **0** |
| **22** | 41   | 42     | 31    |
|      | **1**  | **0**  | **1** |
| **25** | 44   | 45     | 34    |
|      | **0**  | **0**  | **1** |
| **1**  | 20   | 21     | 10    |

|      | **19** | **20** | **9** |
|------|--------|--------|-------|
|      | **1**  | **1**  | **0** |
| **22** | 41   | 42     | 31    |
|      | **1**  | **0**  | **1** |
| **25** | 44   | 45     | 34    |
|      | **0**  | **0**  | **1** |
| **1**  | 20   | 21     | 10    |

*Permute* →

|      | **20** | **19** | **9** |
|------|--------|--------|-------|
|      | **0**  | **1**  | **1** |
| **25** | 45   | 44     | 34    |
|      | **1**  | **1**  | **0** |
| **22** | 42   | 41     | 31    |
|      | **0**  | **0**  | **1** |
| **1**  | 21   | 20     | 10    |

# Implementation

*Input Occup*
{10, 20, 30}

*Output Occup*
{20, 25, 15}

*Sorter*

*Sorter*

{3, 2, 1}

{2, 1, 3}

*Raw Requests*

$$\begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

**Permuted Requests**

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

*Match*

*X Bar*

*X Bar*

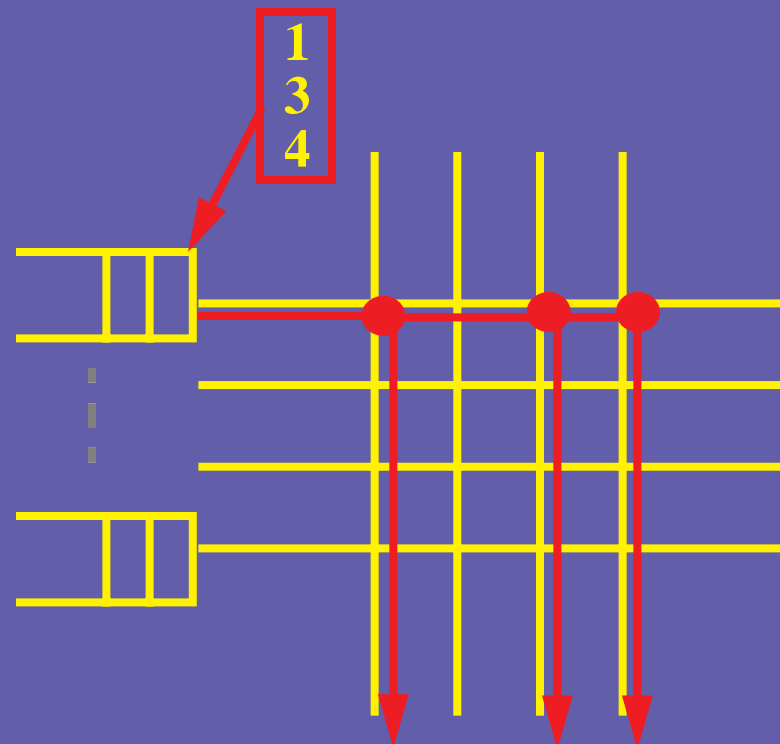*Maxsize Matching*

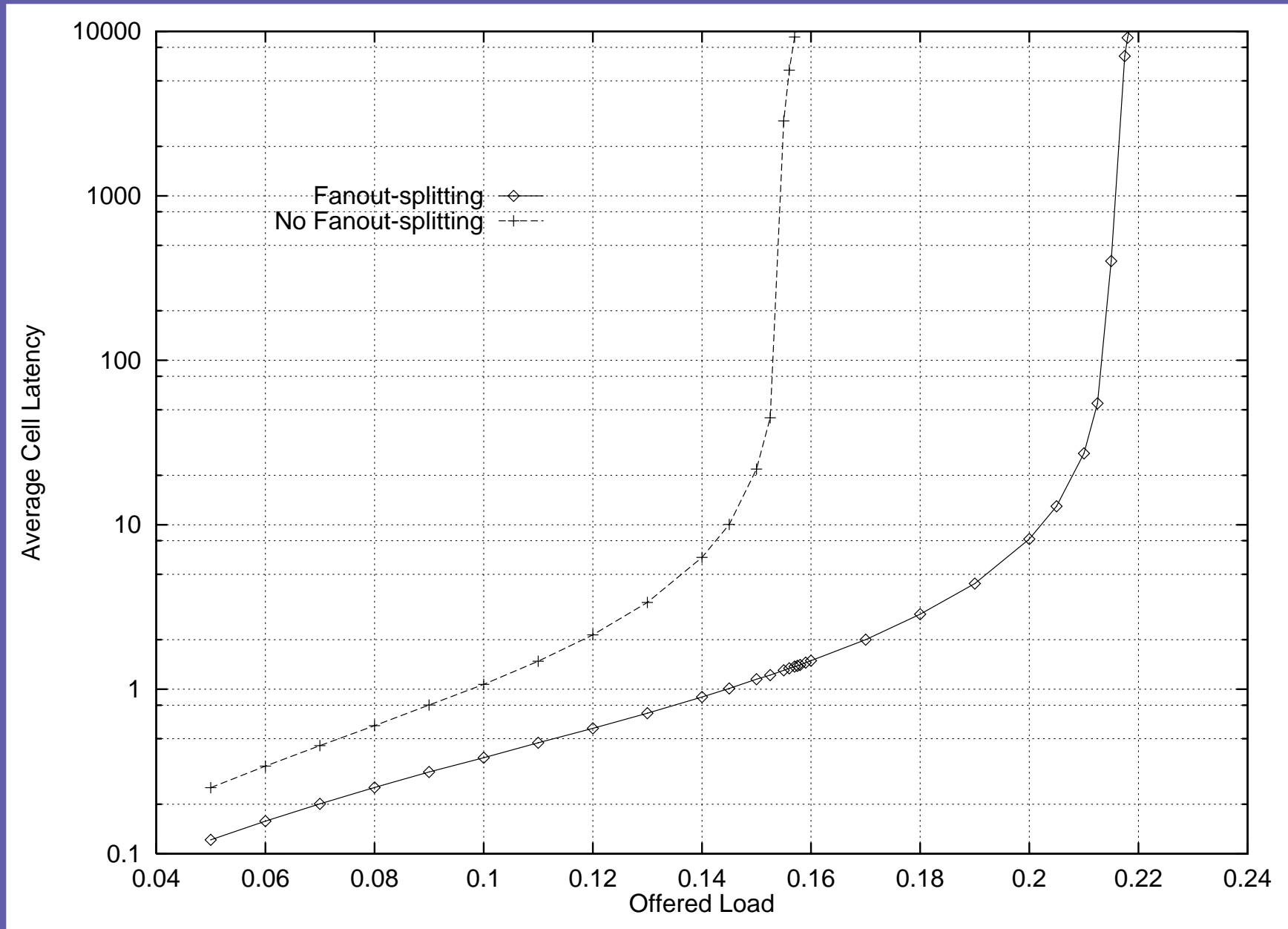*Input permutation*

*Output permutation*

# Multicast Traffic

## Queue Architecture

1. **Making use of the crossbar**

2. **Why treat multicast differently?**

3. **Why maintain a single FIFO queue?**

4. **Fanout-splitting**

# Fanout-Splitting

# Multicast Traffic

## 1. *Residue Concentration*

## 2. *Tetris-based schedulers*
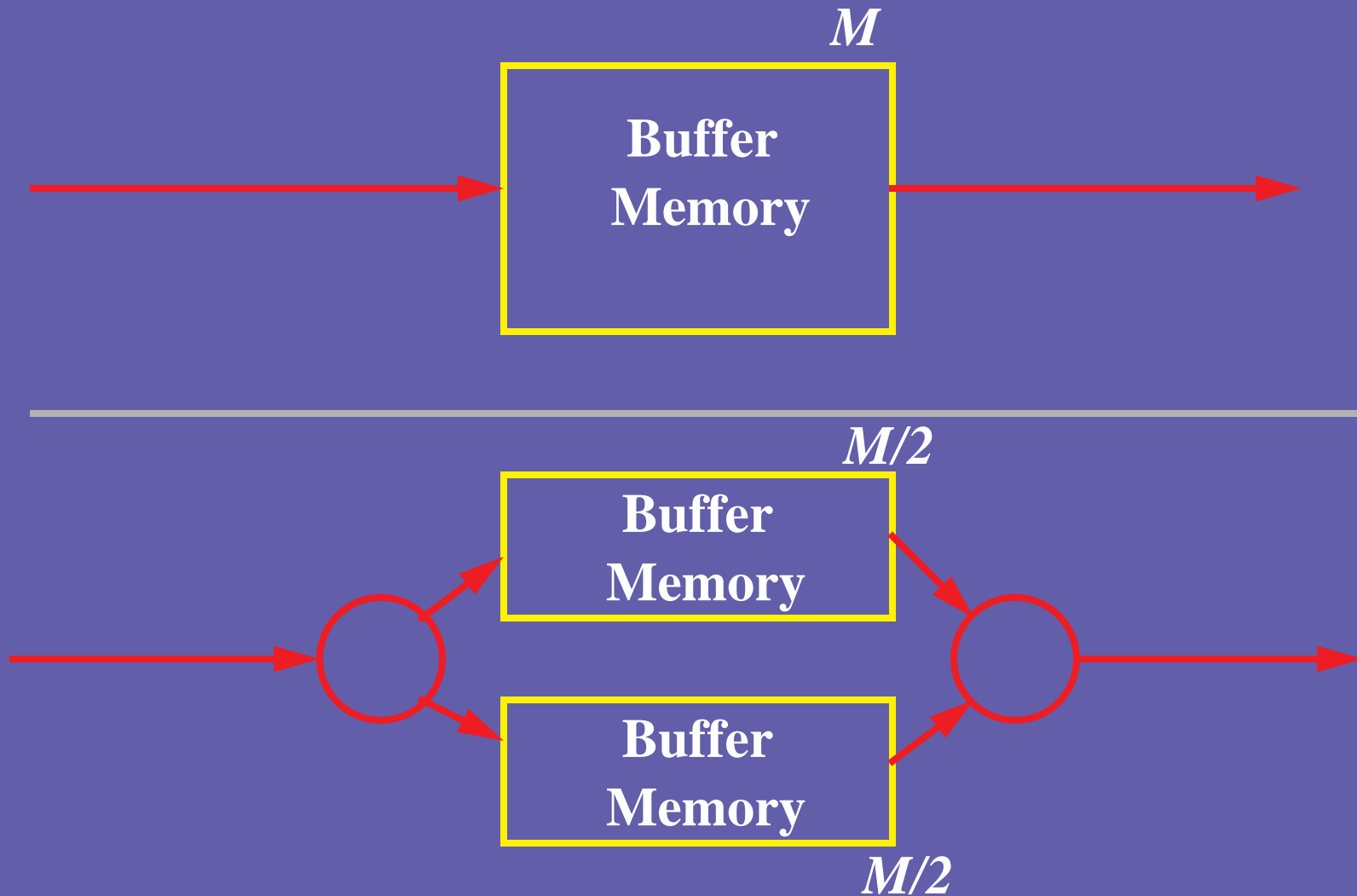
# 1. Accelerating Forwardng Decisions:

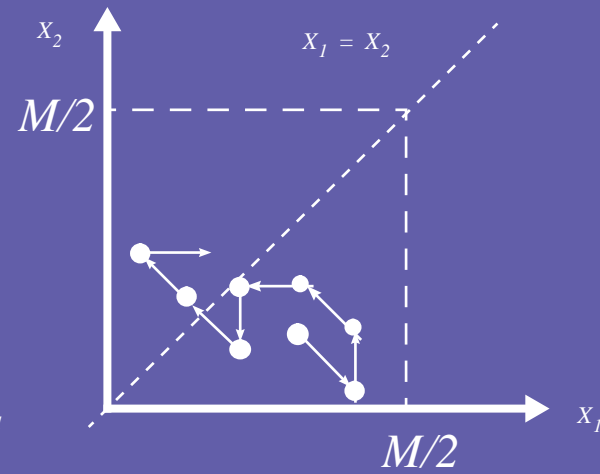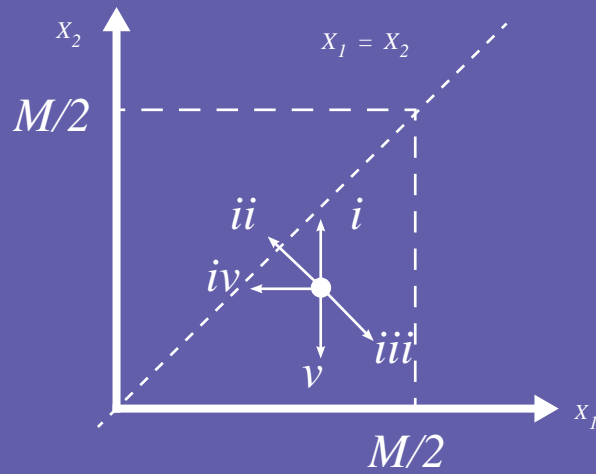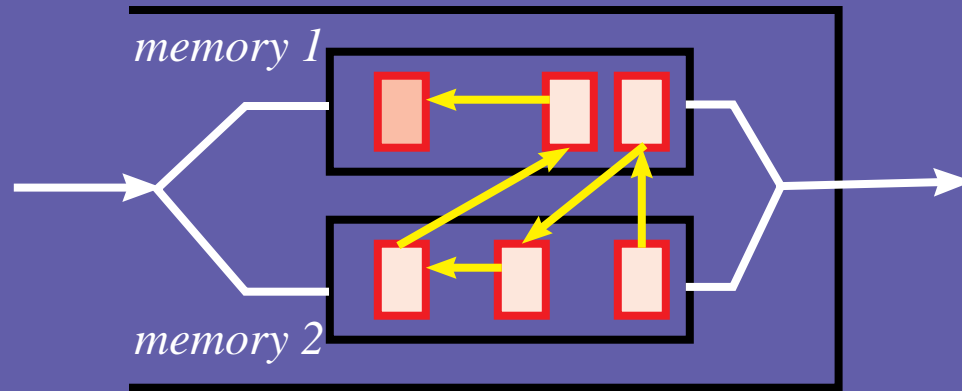- Longest-matching prefixes

# 2. Interconnections: Switched Backplanes

- Input Queueing
  - Theory
  - Unicast
  - Multicast

- Fast Buffering

- Speedup

- The *Tiny Tera* Project
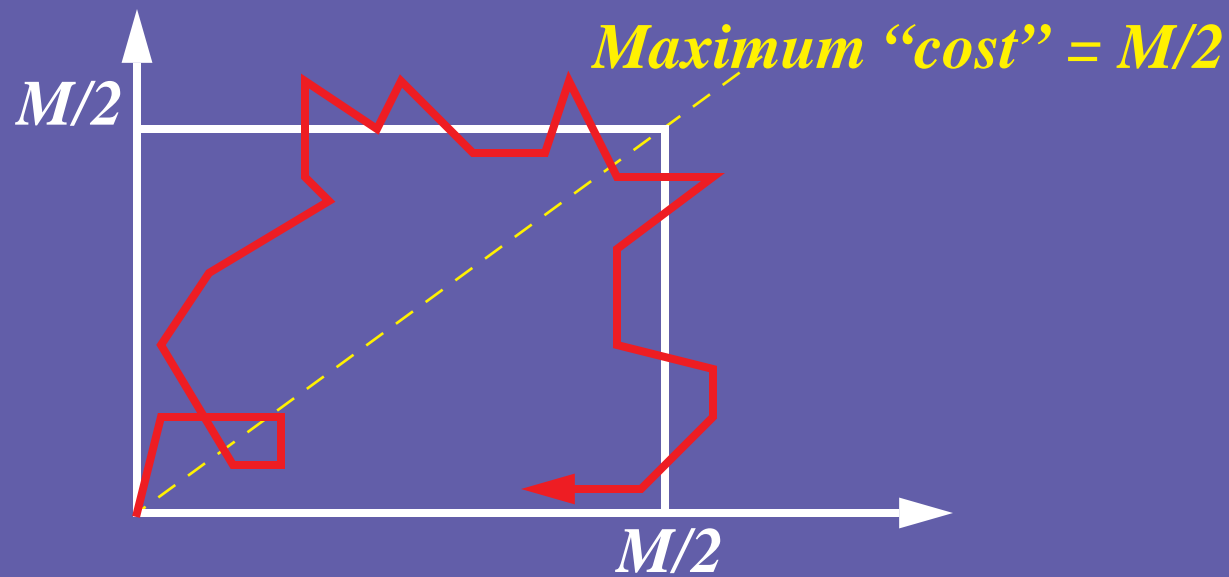
# Fast Buffering
## *Ping-pong Memory*

**M**

Buffer
Memory

**M/2**

Buffer
Memory

Buffer
Memory

**M/2**

*memory 1*

*memory 2*

$X_2$

$X_1 = X_2$

$M/2$

*ii*   *i*

*iv*

*iii*

*v*

$M/2$   $X_1$

$X_2$

$X_1 = X_2$

$M/2$

$M/2$   $X_1$

# Fast Buffering
## *Ping-pong Memory*



Occupancy

M

t

Maximum "cost" = M/2

M/2

M/2

# Fast Buffering
## *Ping-pong Memory*



*single memory: M*

*single memory: M/2*

*ping-pong: (M/2,M/2)*

*In practise, cost <5%*

Log(Overflow Rate)

Buffer size, M

# Some Results
## *Input Queued Switch*

$$Wastage\ Factor,\ \omega(R) \equiv \frac{M(R) - \tilde{M}(R)}{M(R)}$$

- $\omega(R)$ decreases with M

- $\omega(R)$ decreases with burstiness

- $\omega(R)$ decreases with load

- $\omega(R)$ decreases with number of ports

# 1. Accelerating Forwardng Decisions:

- Longest-matching prefixes

# 2. Interconnections: Switched Backplanes

- Input Queueing
  - — Theory
  - — Unicast
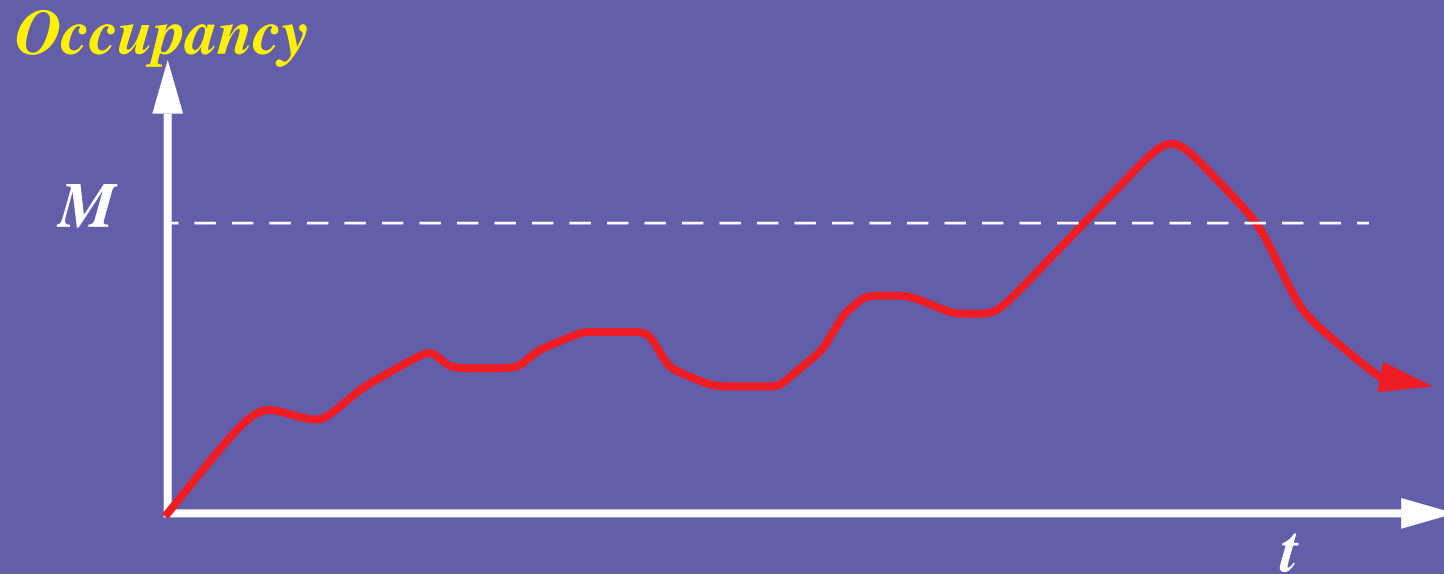  - — Multicast

- Fast Buffering
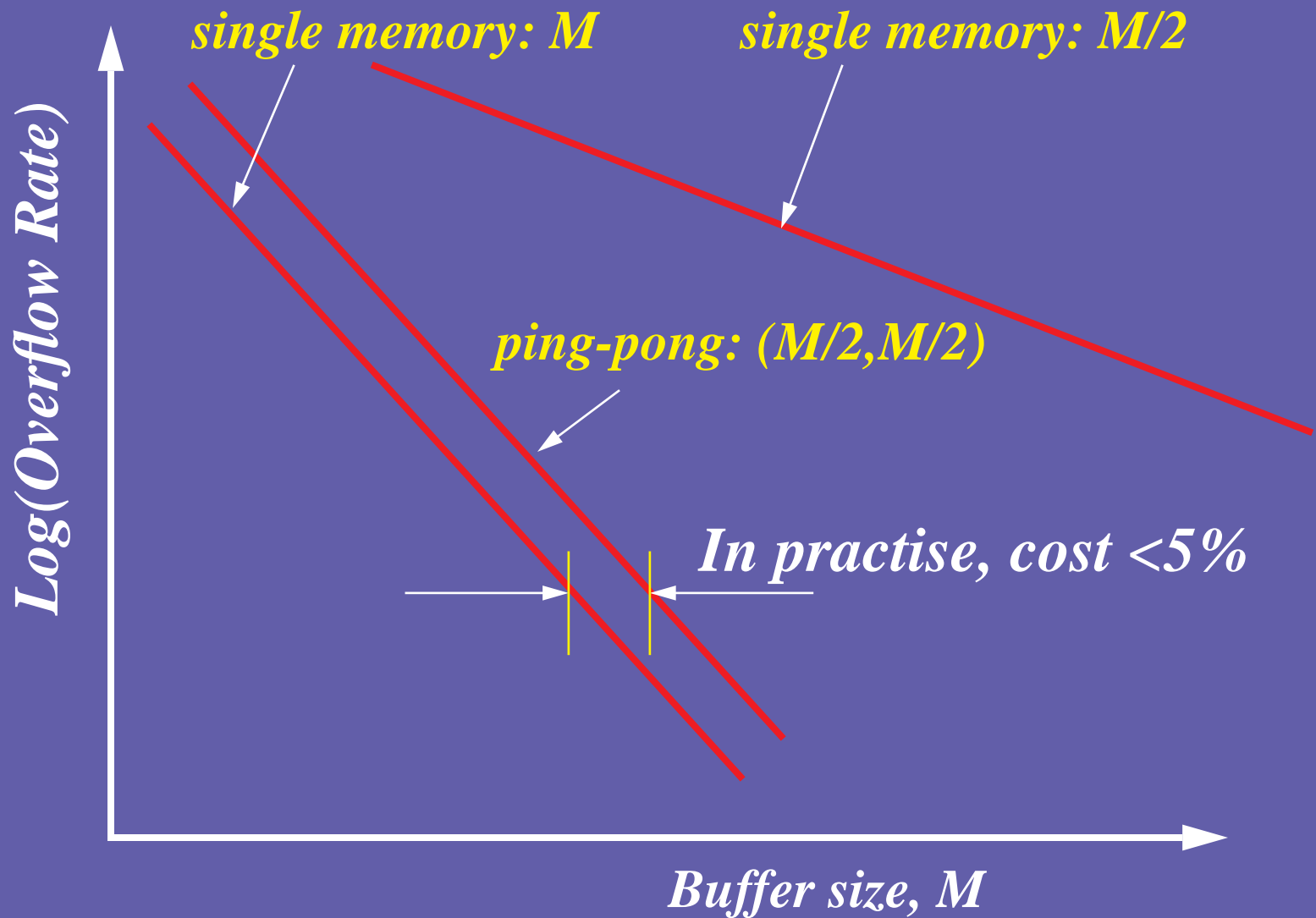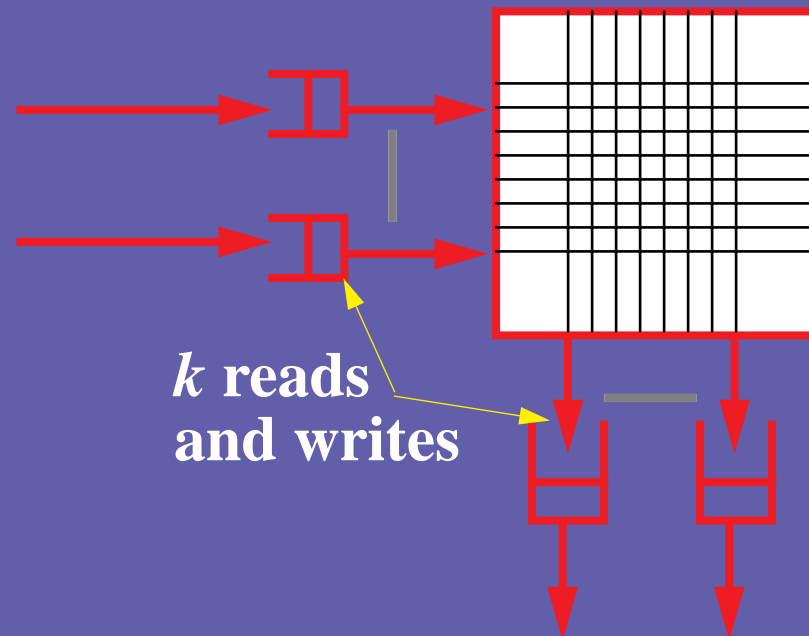
- Speedup

- The *Tiny Tera* Project

# Matching Output Queueing with Input- and Output- Queueing

*How much speedup is enough?*
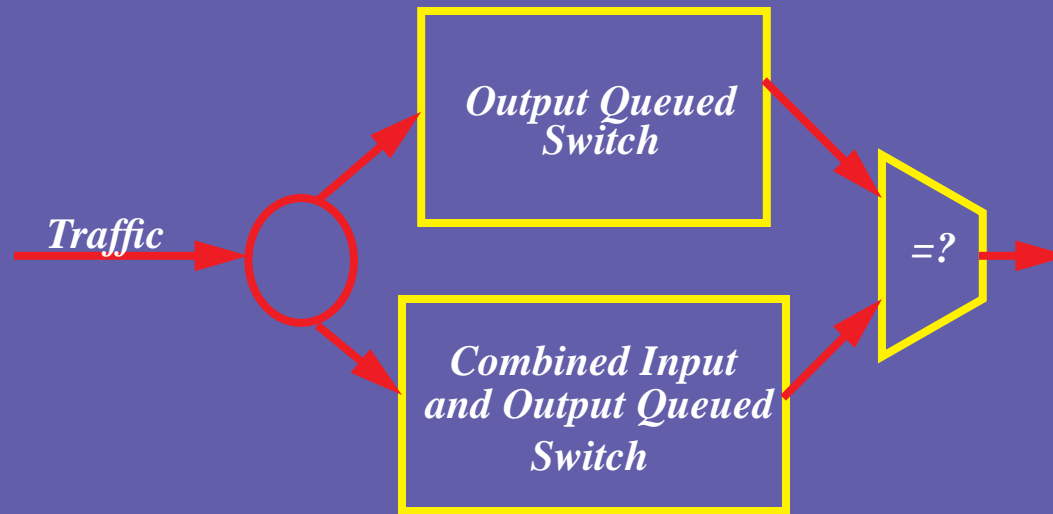
**Combined Input- and Output-Queueing:**



*k* reads and writes

# Matching Output Queueing
# with Input- and Output- Queueing

*How much speedup is enough?*

**Conventional wisdom suggests:**

*A speedup $k = 2 - 4$ leads to high throughput*

# Matching Output Queueing with Input- and Output- Queueing



**Fact**  *To match output queueing, with FIFO input queues:*

$$k = N \quad \text{is necessary and sufficient.}$$

**Fact**  *To match output queueing, with virtual output queues:*

$$k = \left( 2 - \frac{1}{N} \right) \quad \text{is necessary and sufficient.}$$

**1. Accelerating Forwardng Decisions:**

- Longest-matching prefixes

**2. Interconnections: Switched Backplanes**

- Input Queueing
  — Theory
  — Unicast
  — Multicast

- Fast Buffering

- Speedup

- The *Tiny Tera* Project