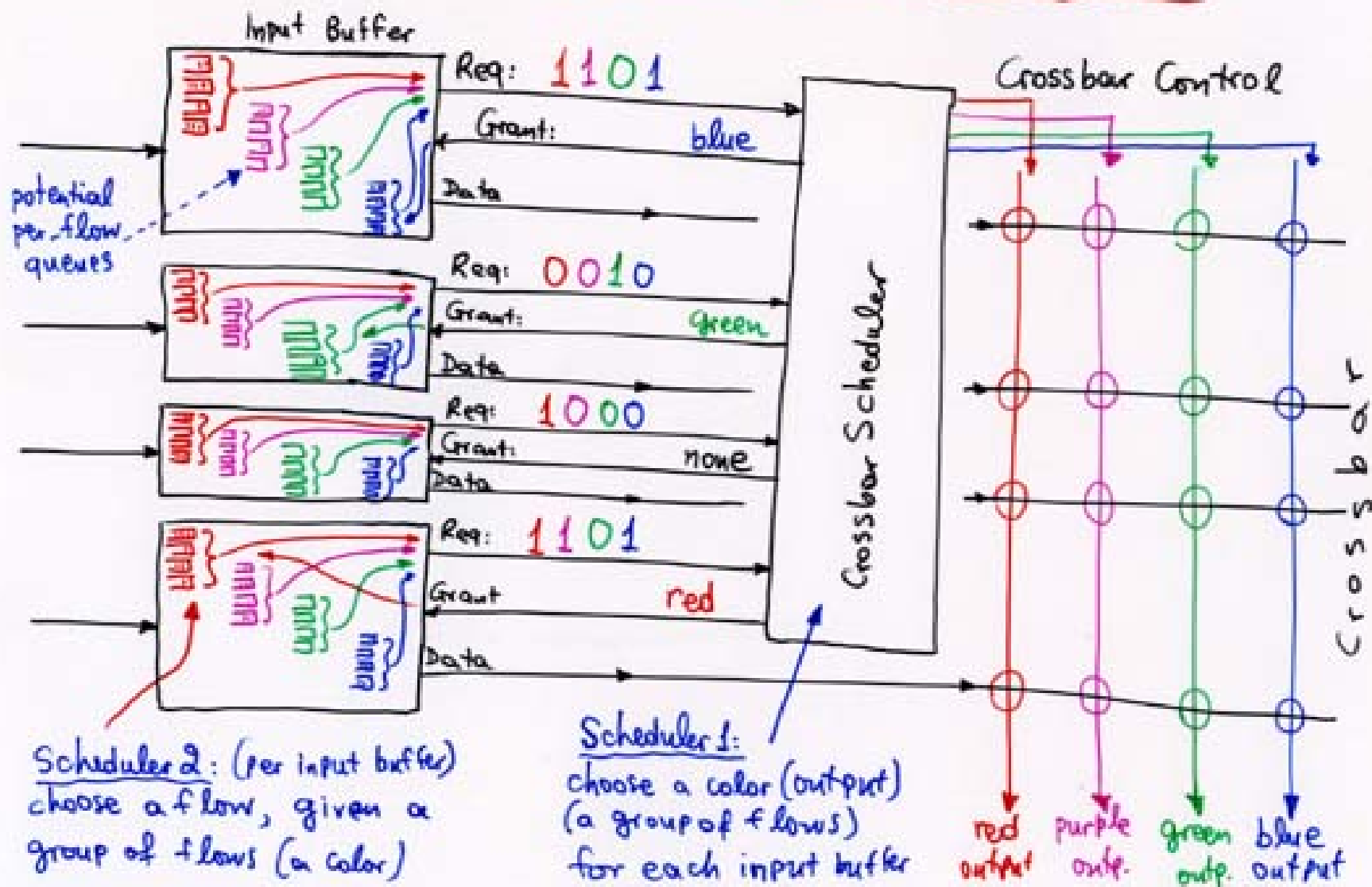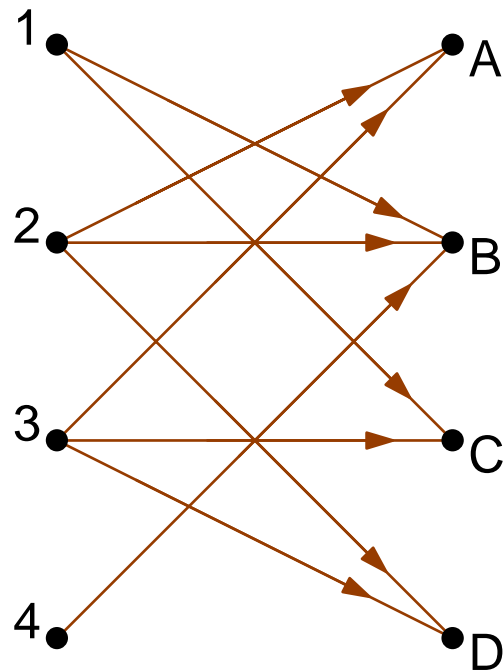# **4.3** Virtual Output Queue (VOQ) Input Q'ing and the Crossbar Scheduling Problem

- Crossbar Switch with one Buffer Memory per Input Line

- Throughput per Buffer Memory = 1 (incoming) + 1 (outgoing)

- Multiple (one per output) Queues per Buffer Memory:

  "Virtual Output Queues – VOQ"

  - $N$ queues per input, $N^2$ queues total, for $N{\times}N$ switch

- Crossbar Scheduling, per cell-time:

  - pairings ("marriages") between inputs and outputs – each input specifies a subset of the outputs that it accepts to be married to

  - interdependent decisions – difficult problem!

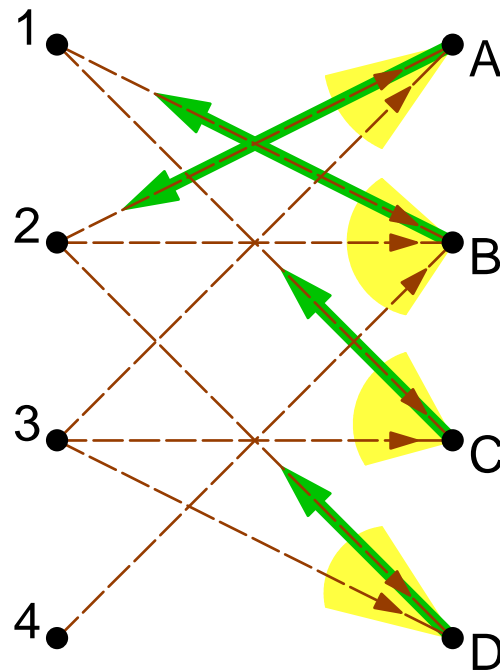Advanced Input Queueing (Buffering) - "Virtual Output" Queues

Input Buffer

Req: 1101
Grant: blue
Data

Req: 0010
Grant: green
Data

Req: 1000
Grant: none
Data

Req: 1101
Grant red
Data

potential per flow queues

Crossbar Scheduler

Crossbar Control

Crossbar

Scheduler 2: (per input buffer) choose a flow, given a group of flows (a color)

Scheduler 1: choose a color (output) (a group of flows) for each input buffer

red output    purple outp.    green outp.    blue output

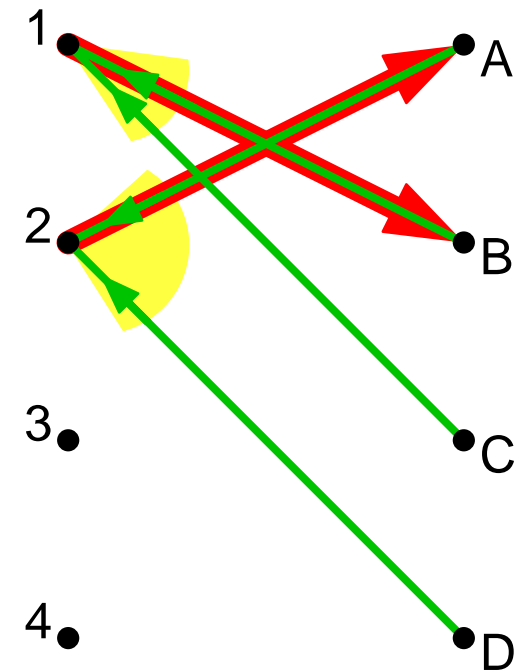# Crossbar Scheduling: Parallel Itarative Matching (PIM)



**Request Phase:**

All inputs send their requests in parallel

**Grant Phase:**
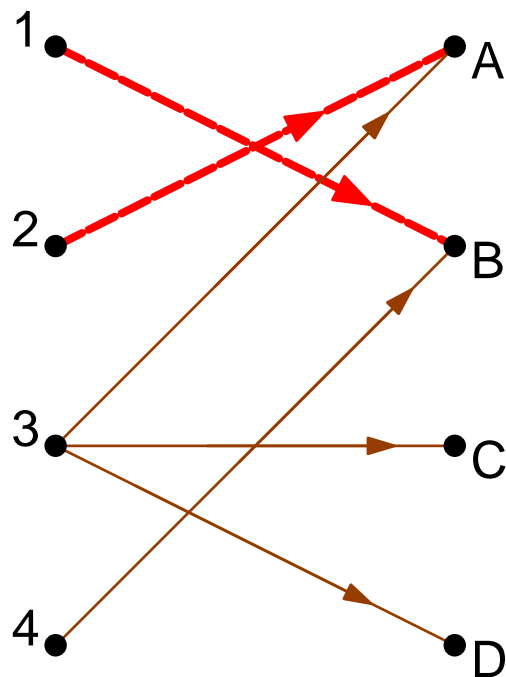
Each output, *independently,* grants to one of the requests that it received
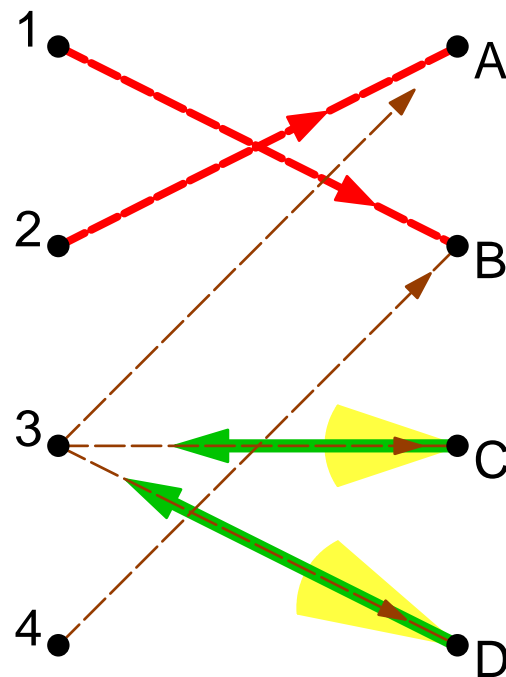
**Accept Phase:**

Each input accepts one of the grants that it received
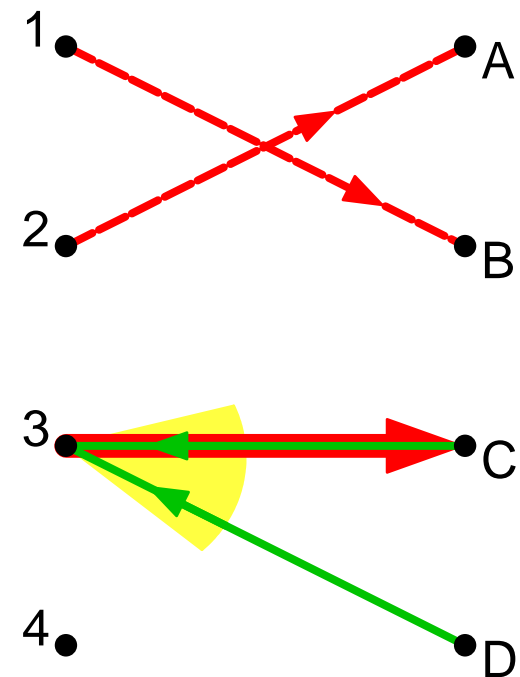
*F i r s t   I t e r a t i o n*

**Request Phase:**
Unmatched inputs
(received no grant)
resend their requests

**Grant Phase:**
Unmatched outputs (rcv'd
no accept) grant to one
of the received requests

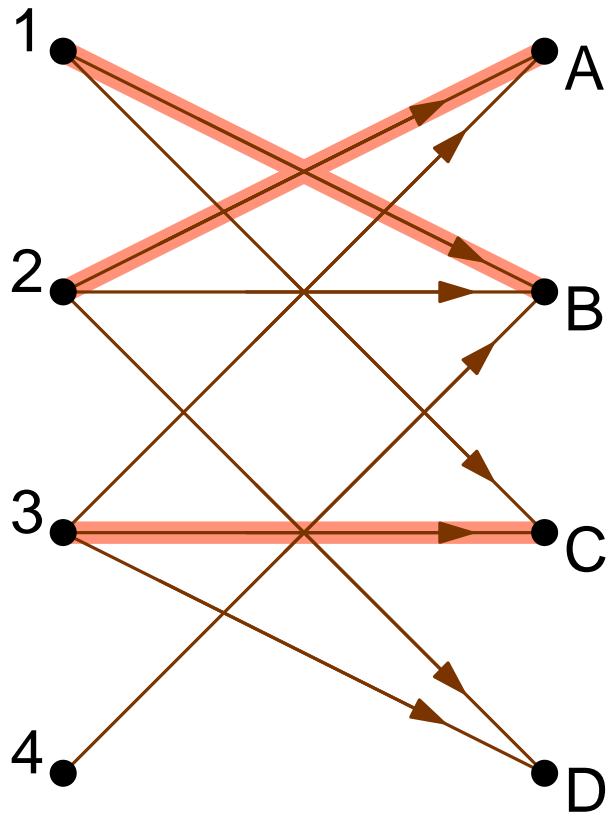**Accept Phase:**
Unmatched inputs
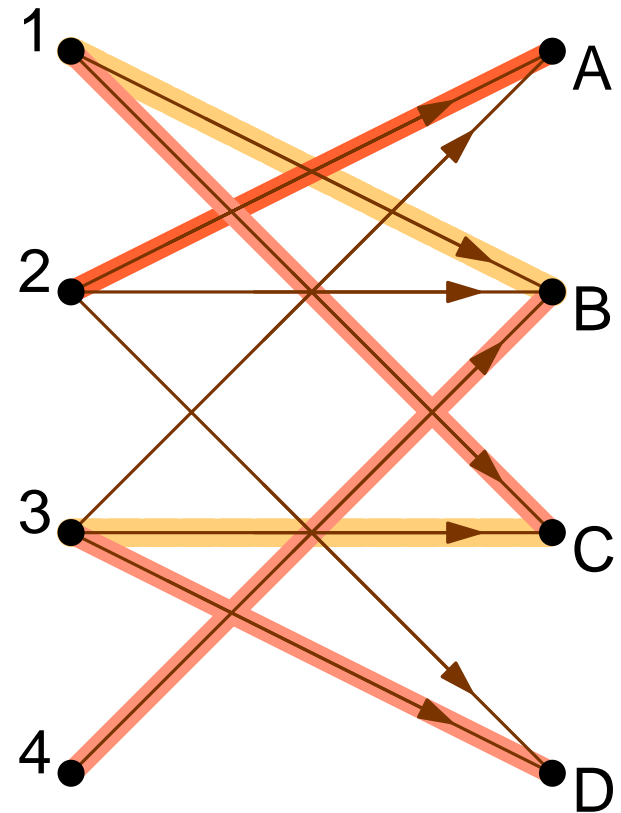accept one of the
received grants

*Second Iteration*

- Original PIM proposal: outputs grant randomly among requesting inputs, inputs accept randomly among granting outputs

## Maximal Matching

Cannot add any new connection
without breaking some
already made connection(s)

## Maximum Matching

Maximum possible number
of connections
for the given request pattern

# Maximum Matching is complex ... and may be Unfair:
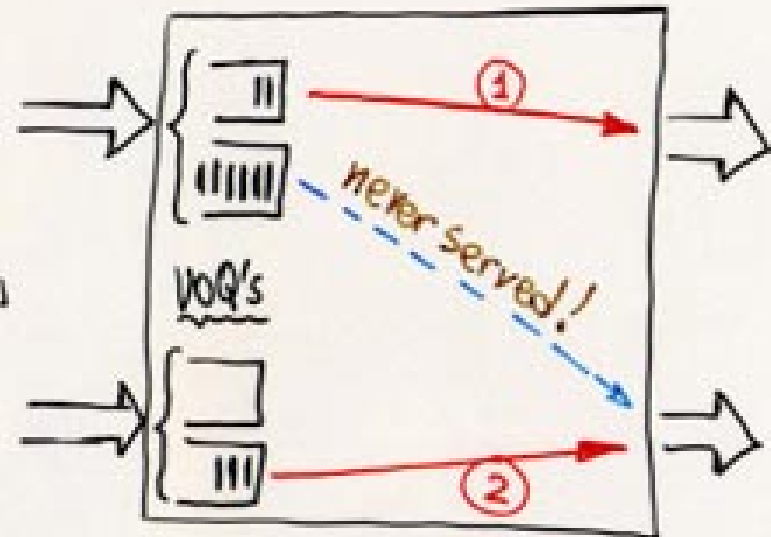
- NxN switch }
- M requests }:

- $O(N \cdot (N+M))$ deterministic algorithm

    Tarjan: Data Structures & Network Algorithms
    SIAM, 1983

- $O(N+M)$ randomized algorithm

    Karp e.a.: ACM STOC, 1990



VOQ's

never served!!

① ②

# Performance of PIM:

- each iteration resolves, on average, $\geq \frac{3}{4}$ of the remaining unresolved requests

$$\Downarrow$$

$$O(\log N) \text{ iterations}$$

- 16×16 switch:

| #iterations | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| % matched | 64-87% | 88-100% | 97-100% | 99.9-100% |

- delay (VOQ w. PIM) $\approx$ 2× delay (Output Queueing)

- how to implement <u>random</u> selection ???

# Fairness Issue:

¼ grant probability per request



¼ accept probability per grant

(<u>weighted</u> random selection ???)

# *iSlip:* Most Popular, Practical Crossbar Scheduler

- Practical variation of PIM

- Widely used in commercial switch products

- Nick McKeown: "The iSLIP Scheduling Algorithm for Input-Queued Switches", IEEE/ACM Tr. on Networking, April 1999

- Performance properties:

  – performs well under uniform and heavy load, when most VOQ's are non-empty and matching almost "rotate" among inputs & outputs

  – adds delay under medium loads, until most VOQ's become non-empty

  – does not perform very well under *"unbalanced"* traffic (each input preferentially sends to one or a few "favored" output(s) of itself)

**iSLIP** → variation of PIM ⟨ fewer iterations / better fairness

Nick McKeown: IEEE/ACM ToN April 1999

• was used in CISCO GSR-12000, Tiny Tera, Abrizio/PMC-Sierra, e.a.

(a) <u>REQUEST</u>: like PIM

(b) <u>GRANT</u>:
- PIM: each output <u>randomly</u> grants to a requesting ⌐input
- <u>iSLIP</u>: grant in a <u>Round-Robin</u> fashion
  to a requesting input;
  top priority = <u>next of</u> : ~~previous grant~~ **NO.!!! – see below**
  <u>previous grant</u> that <u>was</u> <u>accepted</u>!

(c) <u>ACCEPT</u>:
- PIM: each input <u>randomly</u> accepts one of the granting outputs
- <u>iSLIP</u>: accept in <u>Round-Robin</u> priority a granting output,
  with top priority = <u>next of</u> <u>previous accept</u>

• if we were granting round-robin after the previous grant, without regard as to whether the grant was accepted or not, then the grant pointers may get synchronized in a "bad way" and stay that way for a long time, resulting in very poor perf.

• perform well even with a single iteration !
(versus 2 to 4 iterations of PIM)
reason: under heavy load, with ~all inputs requesting ~all outputs, pointers get and stay de-synchronized and scheduler degenerates into time-division multiplexing
⇒ can reach even 100% throughput under uniform load

Example:



in1

in2

out 1

out 2

grant
t=1

t=3

t=2

t=1

t=3

grant t=2

requests
(persistent)

• more iterations improve the delay
• has good fairness properties
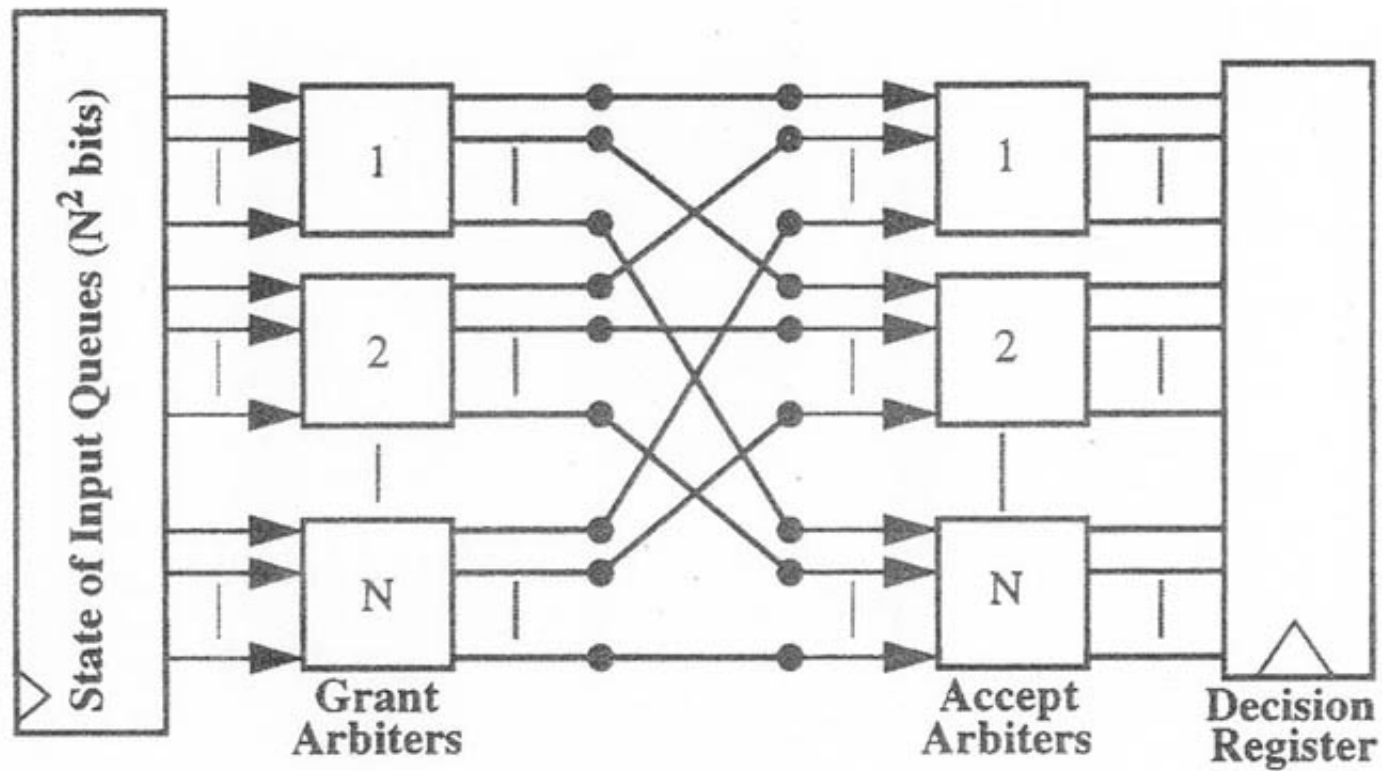• is relatively easy to implement
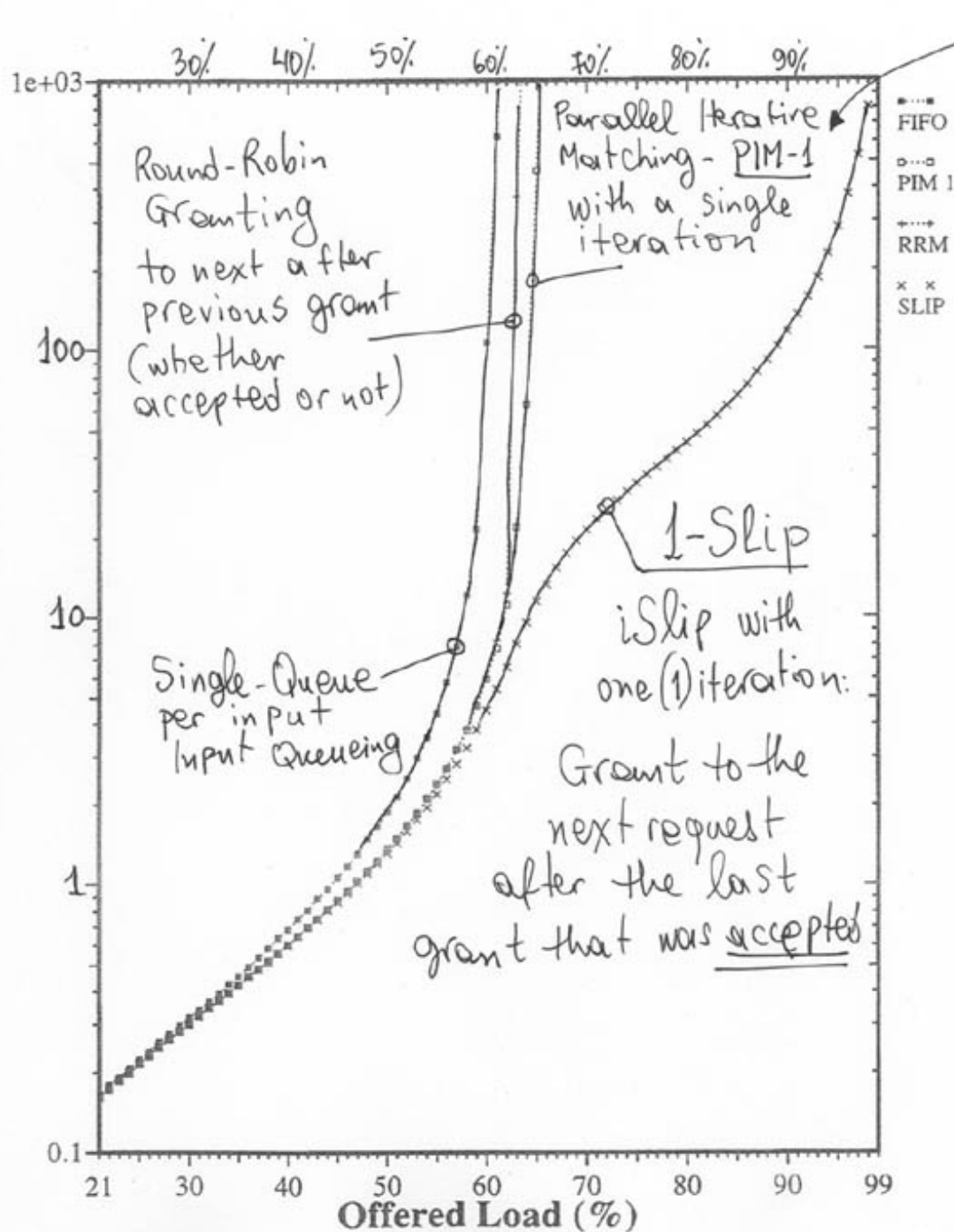
Fig. 21. Interconnection of $2N$ arbiters to implement $i$SLIP for an $N \times N$ switch.

Avg Cell Latency (Cells) vs Offered Load (%) — 16×16 switch simulations

Top axis: 30% 40% 50% 60% 70% 80% 90%

Legend:
- FIFO
- PIM 1
- RRM
- SLIP

Round-Robin Granting to next after previous grant (whether accepted or not)

Parallel Iterative Matching - PIM-1 with a single iteration

Single-Queue per input Input Queueing

1-Slip

iSlip with one (1) iteration:

Grant to the next request after the last grant that was accepted
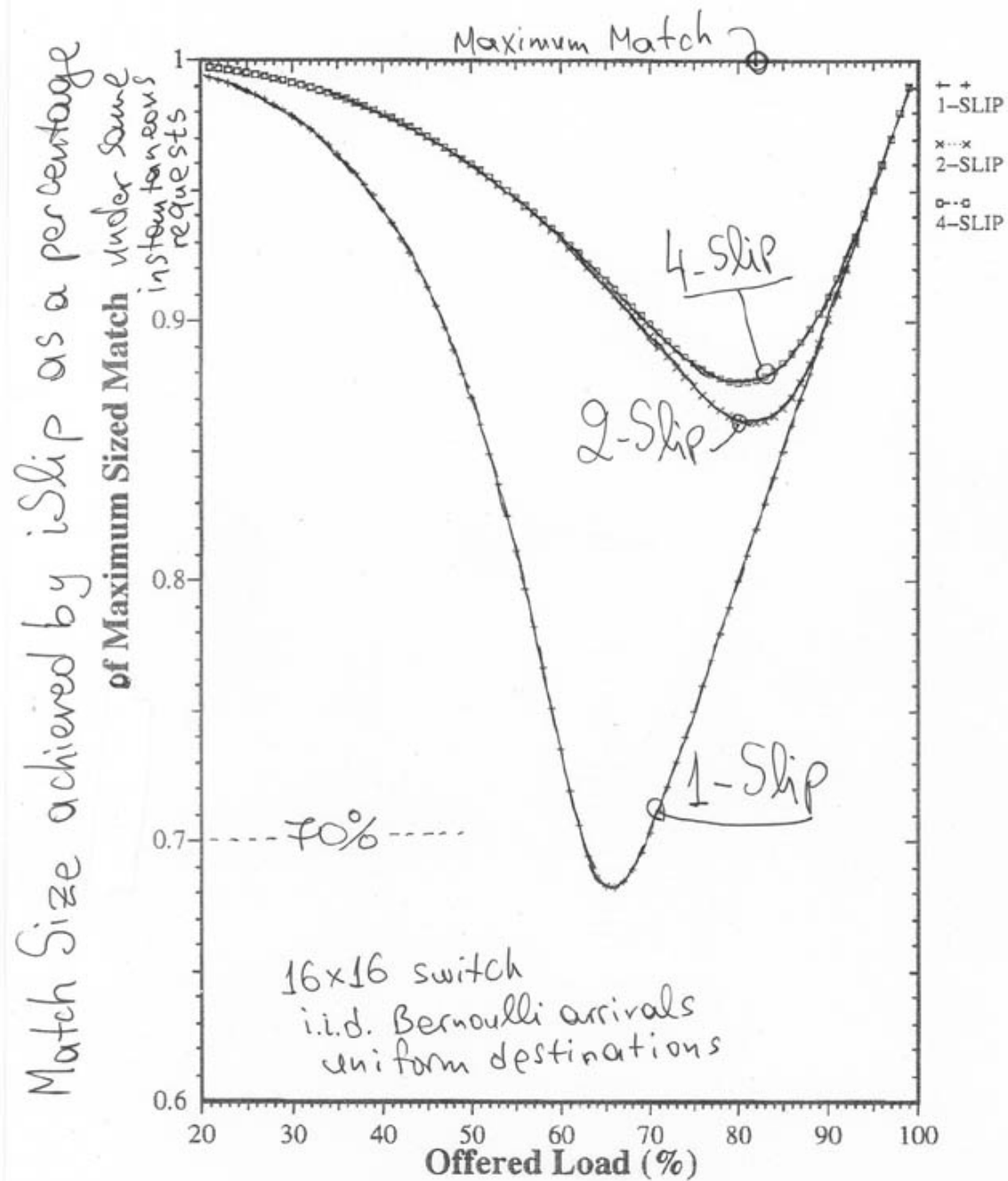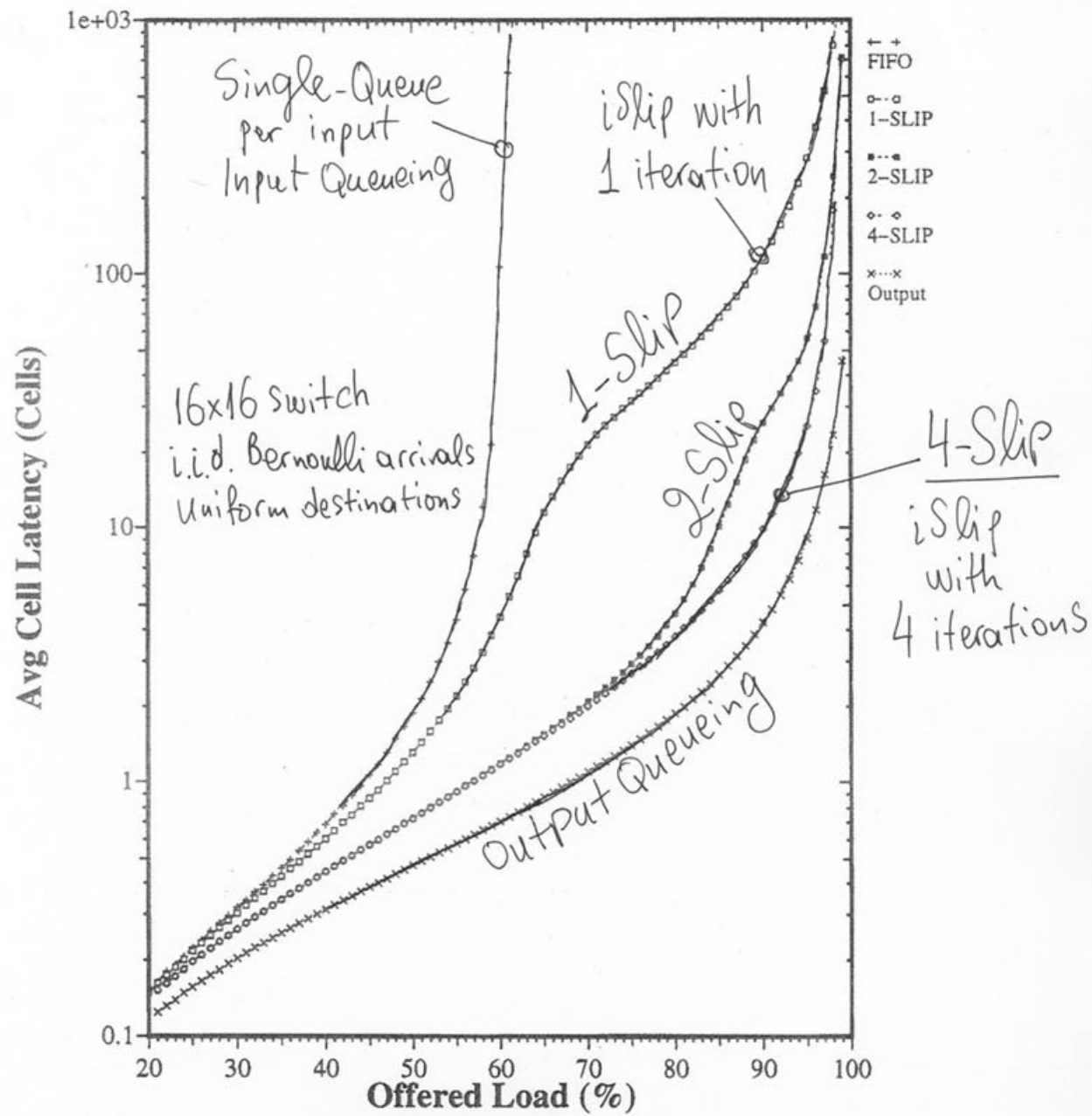
Probability an input remains ungranted

$$= \left(\frac{N-1}{N}\right)^N \text{ for } N \times N \text{ switch}$$

$$N \to \infty \quad \downarrow$$

$$1 - \frac{1}{e} \approx 63\%$$

16x16 switch simulations
i.i.d. Bernoulli arrivals
(non-bursty)
uniformly destined

Y-axis: Avg Cell Latency (Cells) — 1e+03, 100, 10, 1, 0.1

X-axis: Offered Load (%) — 21 30 40 50 60 70 80 90 99

Match Size achieved by iSlip as a percentage of Maximum Sized Match under same instantaneous requests

Maximum Match

1-SLIP

2-SLIP

4-SLIP

4-Slip

2-Slip

1-Slip

70%

16x16 switch
i.i.d. Bernoulli arrivals
uniform destinations

Offered Load (%)

13

14

15

16

# Pipelined Scheduling for Cell-Based Crossbars



Requests In

Grants Out

1 cell time

$Q_{ij}(0) \rightarrow$ Scheduler Instance 0 $G_{ij}(0) \rightarrow$ Cell Transfer

$-1$

$+1$

$Q_{ij}(1) \rightarrow$ Scheduler Instance #1 $G_{ij}(1) \rightarrow$ Cell Xfer

1 cell time

$-1$

$+1$

$Q_{ij}(t) \rightarrow$ Sch. Inst. (t) $G_{ij}(t) \rightarrow$ Cell Xfer

1 cell time

$-1$

$+1$

$Q_{ij}(k-1) \rightarrow$ $G_{ij}(k-1) \rightarrow$

$+1$ $+1$

$-1$

$Q_{ij}(k) \rightarrow$ $G_{ij}(k) \rightarrow$

$+1$ $+1$

$-1$

$Q_{ij}(k+1) \rightarrow$

$+1$ $+1$ $+1$

Cell Arrivals

$Q_{ij}$ = Occupancy (# of cells) of
Queue i→j minus # of pending
requests issued on behalf
of that queue

time

- Cell Arrivals increment $Q_{ij}$
- When $Q_{ij} \geq 1 \Rightarrow$ a Request $R_{ij}$ is issued
- Requests issued decrement $Q_{ij}$ (pending decision on cell transfer)
- When Pending Decision is resolved:
    - Successful Grants leave $Q_{ij}$ as is (already decremented)
    - Failed Grants re-increment $Q_{ij}$ (restore unaccepted request)

<u>Assumption</u> for this scheme to work: <u>fixed-size cells</u>
$\Rightarrow$ requests for cells and actual cell transfers are <u>interchangeable</u> with each other: if the "first" request, on behalf of the "first" cell in the queue is not granted, then the "second" request (issued on behalf of the "second" cell) will result in transfering the "first" cell, if granted.

<u>Ref:</u> Oki, Rojas-Cessa, Chao: "A Pipeline-Based Approach for Maximal-Sized Matching Scheduling in Input-Buffered Switches" IEEE Communications Letters, June 2001 — http://acts.poly.edu/~chao/publications.html