# ECE 6115 / CS 8803 - ICN

## Interconnection Networks for High Performance Systems
### Spring 2020

# FLOW-CONTROL

**Tushar Krishna**

Assistant Professor

School of Electrical and Computer Engineering

Georgia Institute of Technology

tushar@ece.gatech.edu

# NETWORK ARCHITECTURE

- **Topology**
  - How to connect the nodes
  - ~Road Network

- **Routing**
  - Which path should a message take
  - ~Series of road segments from source to destination

- **Flow Control**
  - When does the message have to stop/proceed
  - ~Traffic signals at end of each road segment
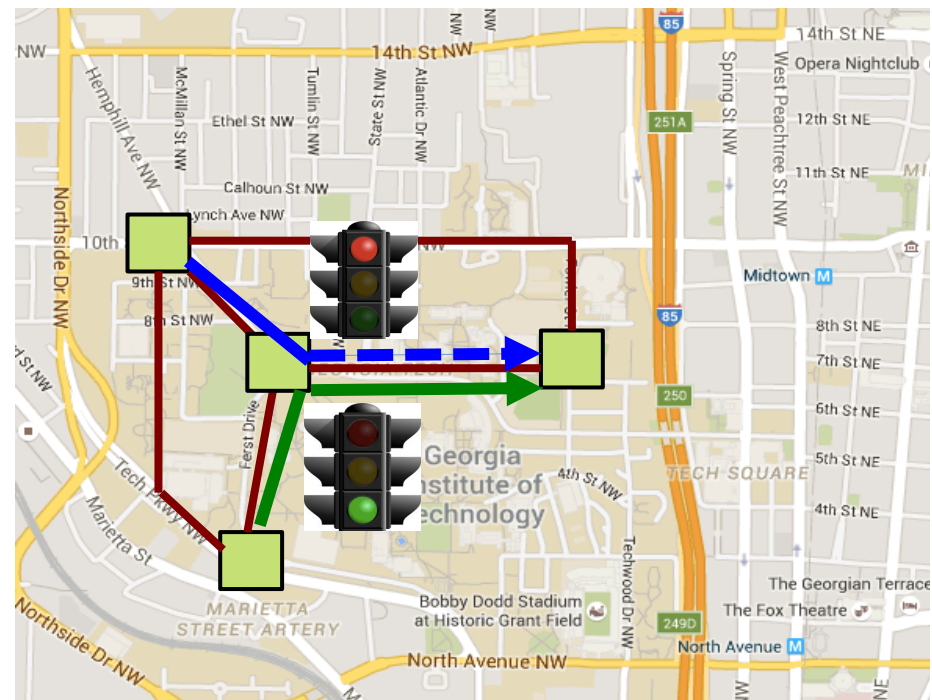
- **Router Microarchitecture**
  - How to build the routers
  - ~Design of traffic intersection (number of lanes, algorithm for turning red/green)

# FLOW CONTROL

Once the topology and route are fixed, flow control determines the **allocation of network resources** (channel bandwidth, buffer capacity, and control state) to packets as they traverse the network
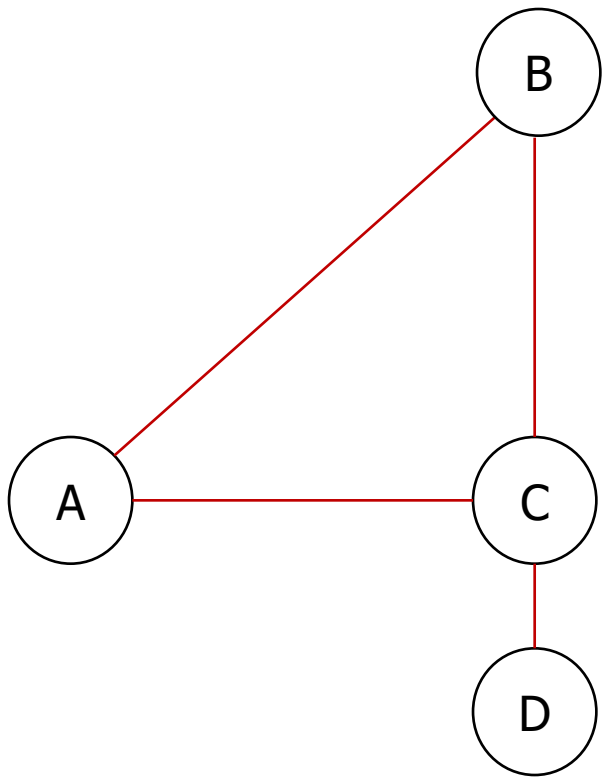
*== resolution of contention between packets requesting the same resource*

~Traffic Signals / Stop signs at end of each road segment

# WHY FLOW CONTROL MATTERS?

Flow control can single-handedly determine performance, however efficient the topology or routing algorithm might be
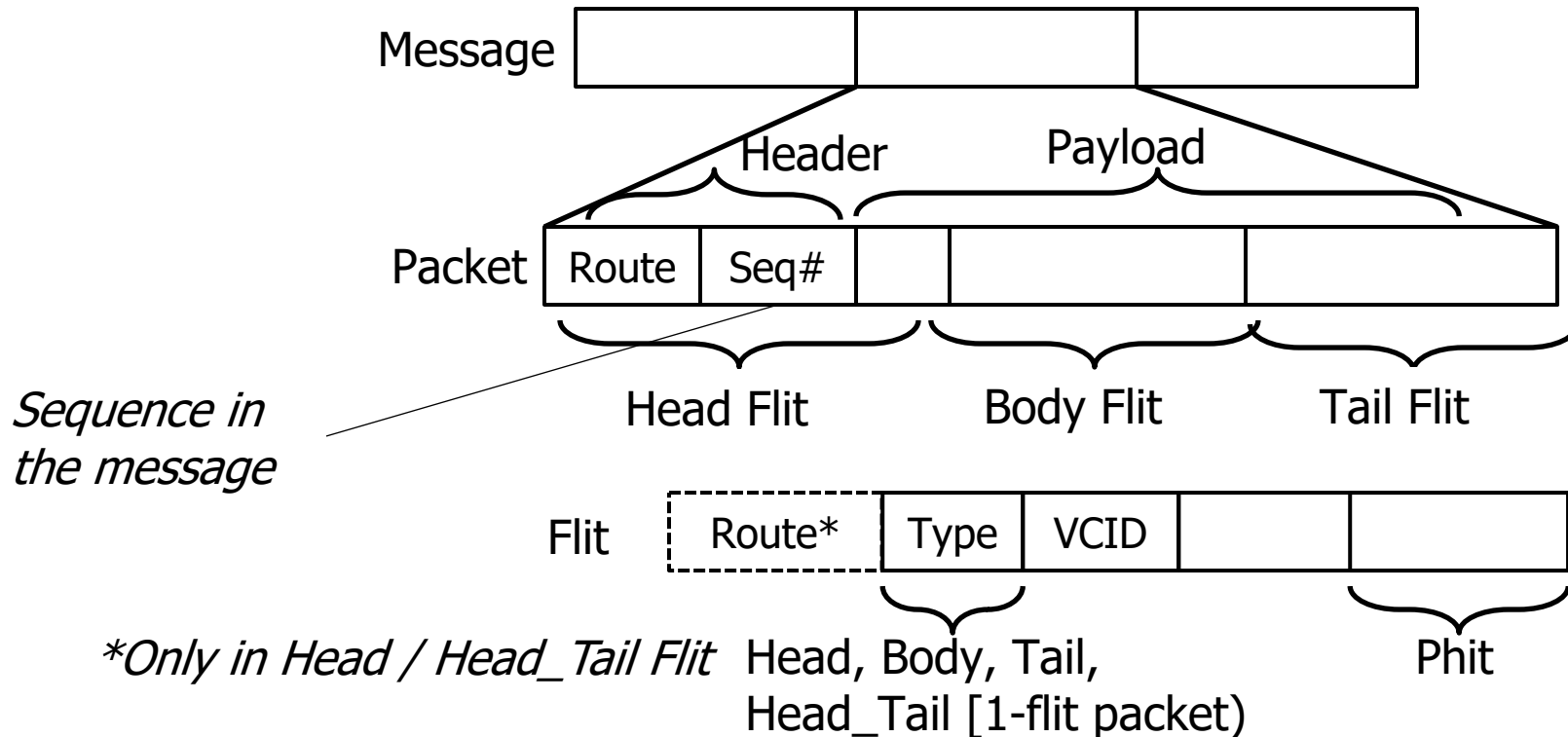


|  | Latency (hops) (A→B) | Throughput (msg/cycle) (A→B) |
|---|---|---|
| **Topology** | 1 | 1 |
| **Routing (XY)** | 2 | 1 |
| **Flow Control** $3$ $(R_A +) L_{AC} + R_C + L_{CB} (+ R_B)$ | | |
| Case I: One buffer at C | | 1/2 |
| Case II: D→B msgs | | 1/5 |

*Case II: D sends 4 messages, 1 cycle break, 4 messages, 1 cycle break...C prioritizes straight over turning traffic*

*Suppose Router Delay = 1, Link Delay = 1*

# ALLOCATION GRANULARITY: MESSAGES, PACKETS, AND FLITS

Message

Header                Payload

Packet | Route | Seq# | | | |

*Sequence in the message*

Head Flit          Body Flit          Tail Flit

Flit | Route* | Type | VCID | | |

*Only in Head / Head_Tail Flit*   Head, Body, Tail, Head_Tail [1-flit packet)
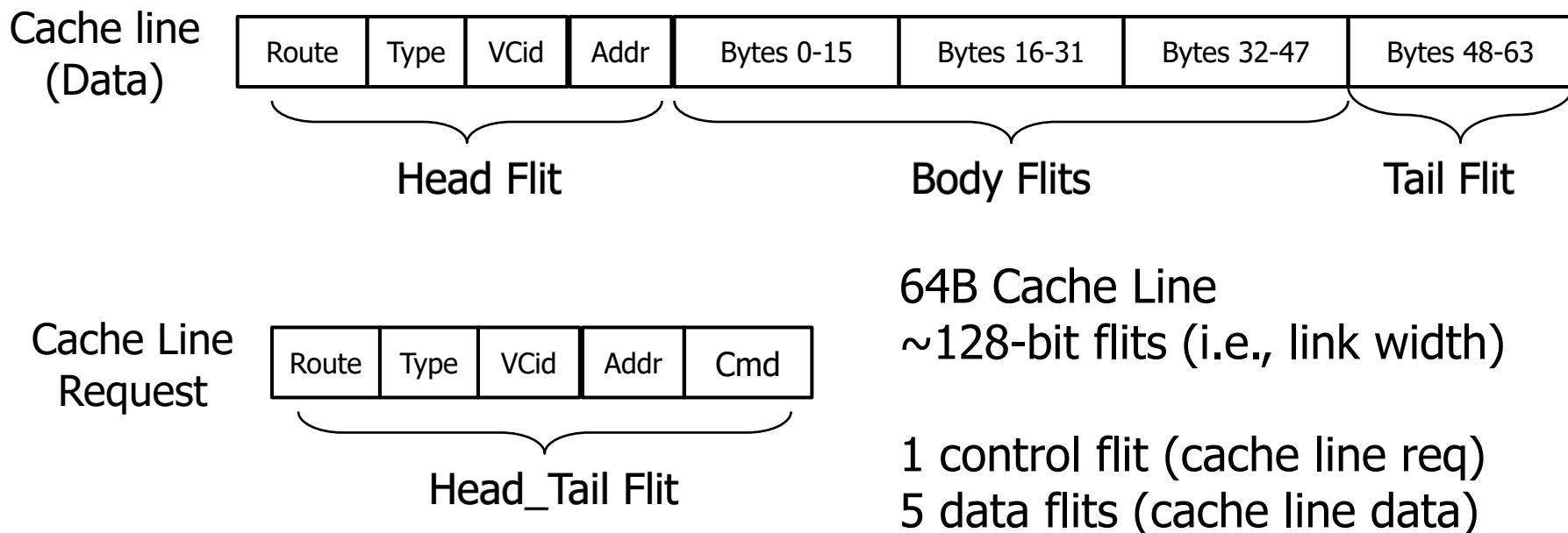
Phit

**Off-chip (SANs)**
Messages could be B/KB/MB of data
Flits have to be sent serially as multiple phits (limited by **pins**)

**On-chip (NoC)**
Message = Packet
Flit = Phit (**abundant on-chip wires**)

# PACKET SIZES IN NOCS

Cache line (Data)

| Route | Type | VCid | Addr | Bytes 0-15 | Bytes 16-31 | Bytes 32-47 | Bytes 48-63 |
|-------|------|------|------|------------|-------------|-------------|-------------|

Head Flit          Body Flits          Tail Flit

Cache Line Request

| Route | Type | VCid | Addr | Cmd |
|-------|------|------|------|-----|

Head_Tail Flit

64B Cache Line
~128-bit flits (i.e., link width)

1 control flit (cache line req)
5 data flits (cache line data)

## All flits of a packet take same route and have the same VCid

# FLOW CONTROL BASED ON ALLOCATION GRANULARITY

- Message-based Flow Control
  - E.g., Circuit Switching

- Packet-based Flow Control
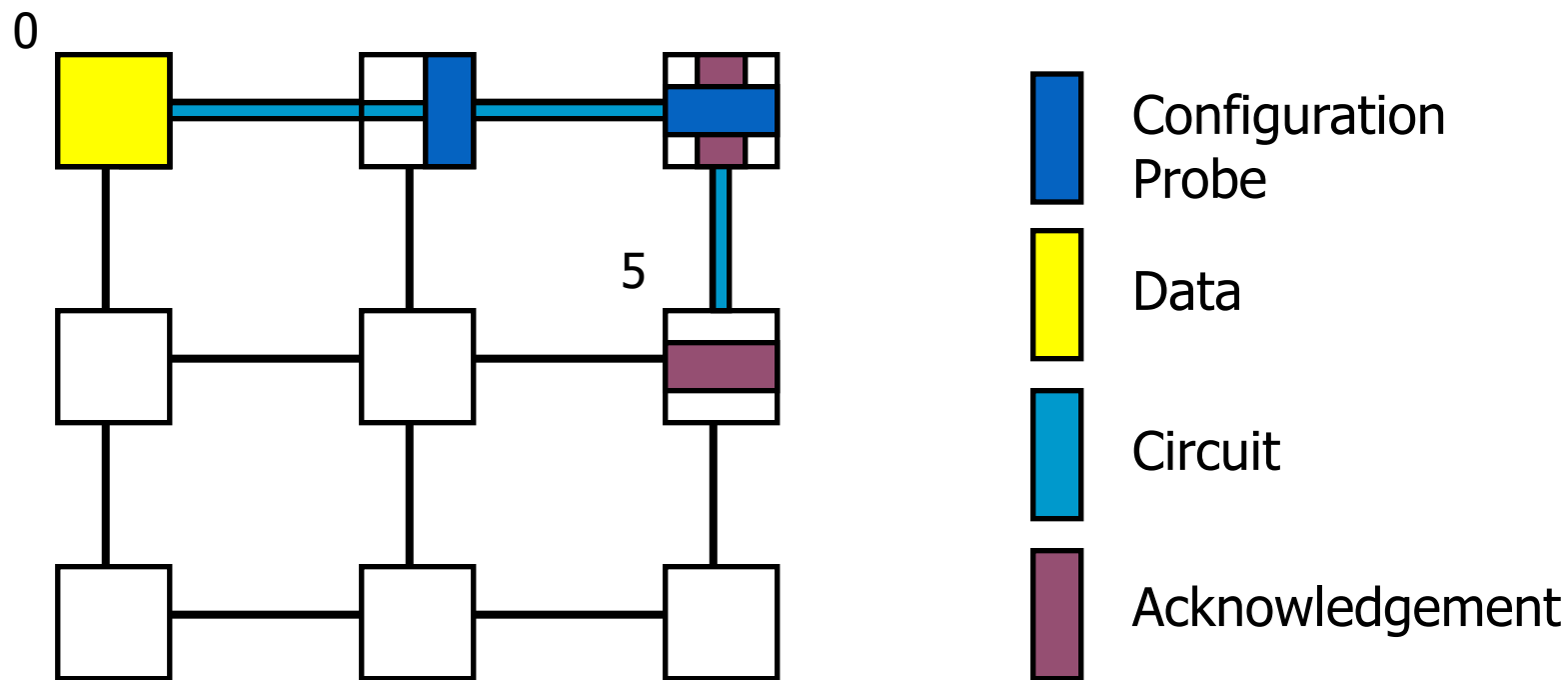  - E.g., Store and Forward, Virtual Cut-Through

- Flit-based Flow Control
  - E.g., Wormhole, Virtual Channel

# MESSAGE-BASED FLOW CONTROL

- Coarsest Granularity

- Circuit-switching
  - **Setup entire path before sending message**
    - Reserve all channels from source to destination using a setup probe
  - **Once setup complete, send Data through the channels**
    - Buffers not needed at routers as no contention
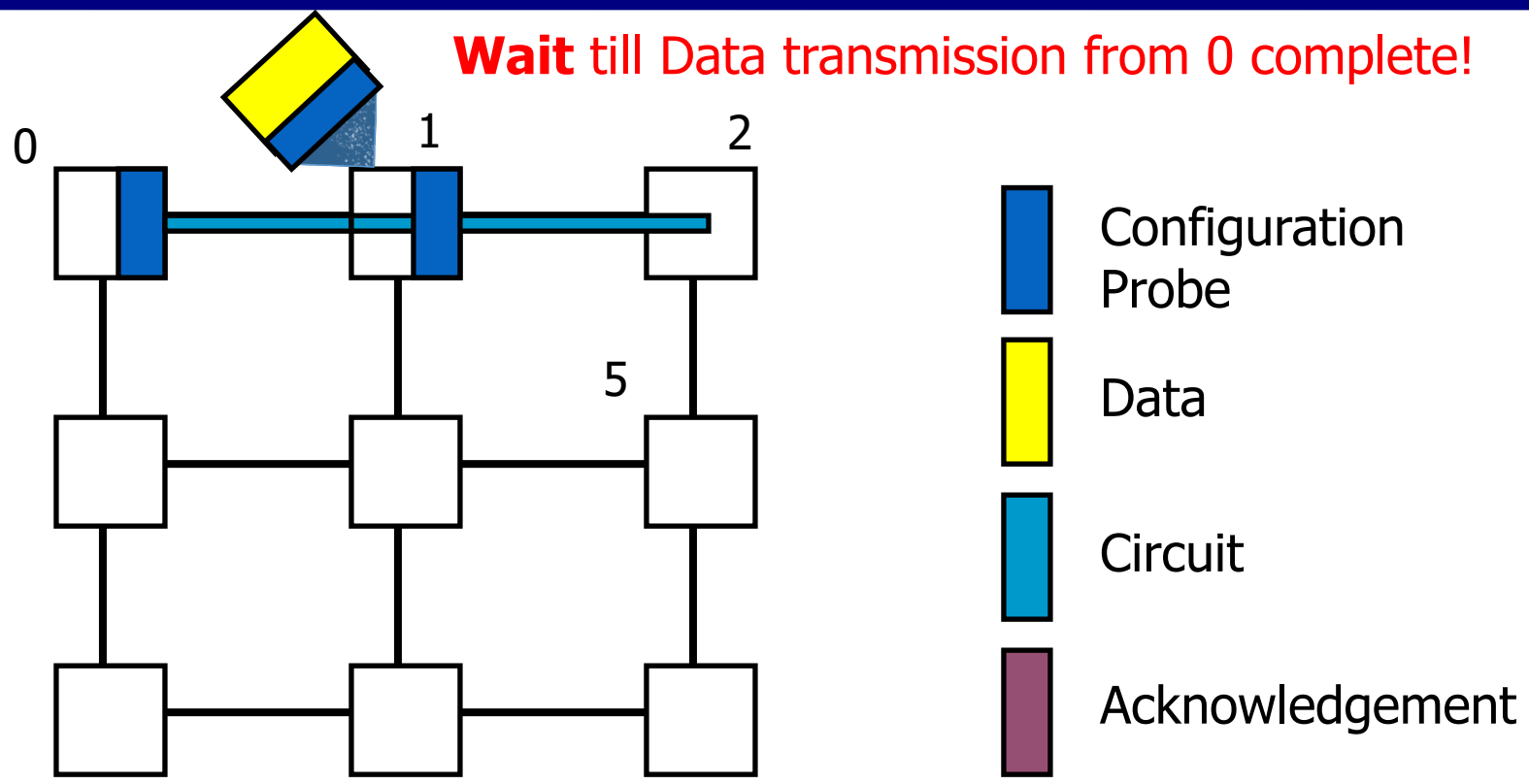  - **Tear down the circuit once transmission complete**

# CIRCUIT SWITCHING EXAMPLE



**Configuration Probe** (blue)

**Data** (yellow)

**Circuit** (cyan)

**Acknowledgement** (maroon)

- Significant latency overhead prior to data transfer
  - Data transfer does not pay per-hop overhead for buffering, routing, and allocation

# HANDLING CONTENTION

**Wait** till Data transmission from 0 complete!



0   1   2

5

- Configuration Probe
- Data
- Circuit
- Acknowledgement

- When there is contention
  - Significant wait time
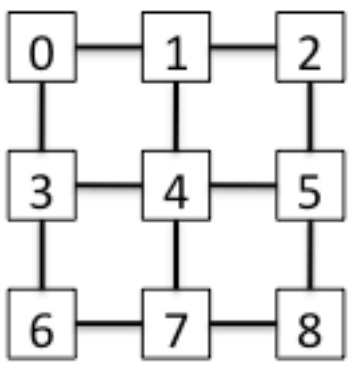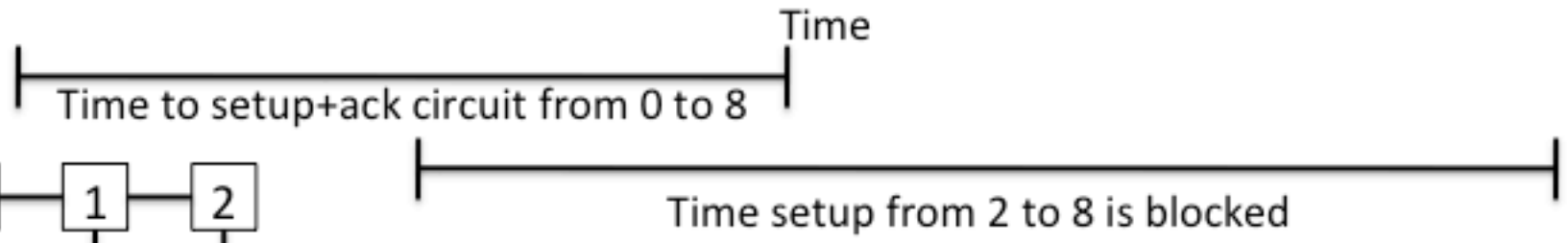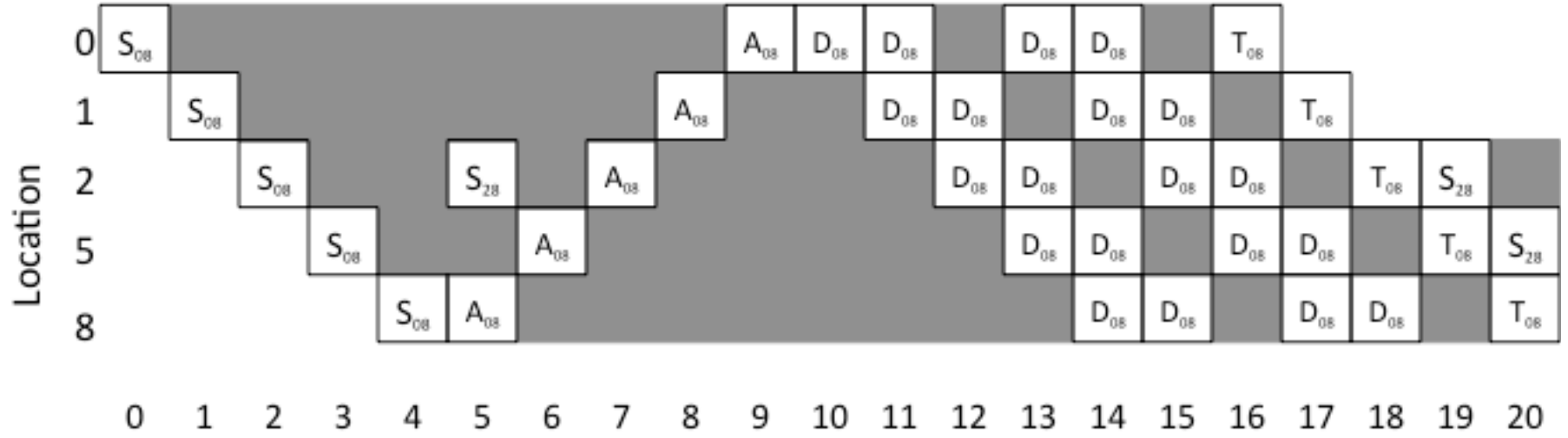  - Message from 1 → 2 must wait

# CHALLENGES WITH CIRCUIT-SWITCHING

- Loss in bandwidth (throughput)
  - Throughput can suffer due to **setup** and **transfer** time for circuits
    - Links are idle until setup is complete
    - No other message can use links until transfer is complete

- Latency overhead in setup if the amount of data being transferred is small

# CIRCUIT-SWITCHING IN NOCS?

- Cache Line = 64B
  - Suppose
    - Channel Width = 128b => 64x8/128 = 4 chunks
    - 3-hop traversal with 1-cycle per hop
  - Setup = 3 cycles
  - ACK = 3 cycles
  - Data Transfer Time = 3 (for first chunk) + 3 (remaining chunks) = 6 cycles
  - Total Time = 12 cycles
    - Half of this went in circuit setup!

- Hybrid Circuit-Packet Switching
  - "Jerger et. al, "Circuit Switched Coherence", NOCS 2008

# TIME-SPACE DIAGRAM: CIRCUIT SWITCHING



Time to setup+ack circuit from 0 to 8

Time setup from 2 to 8 is blocked
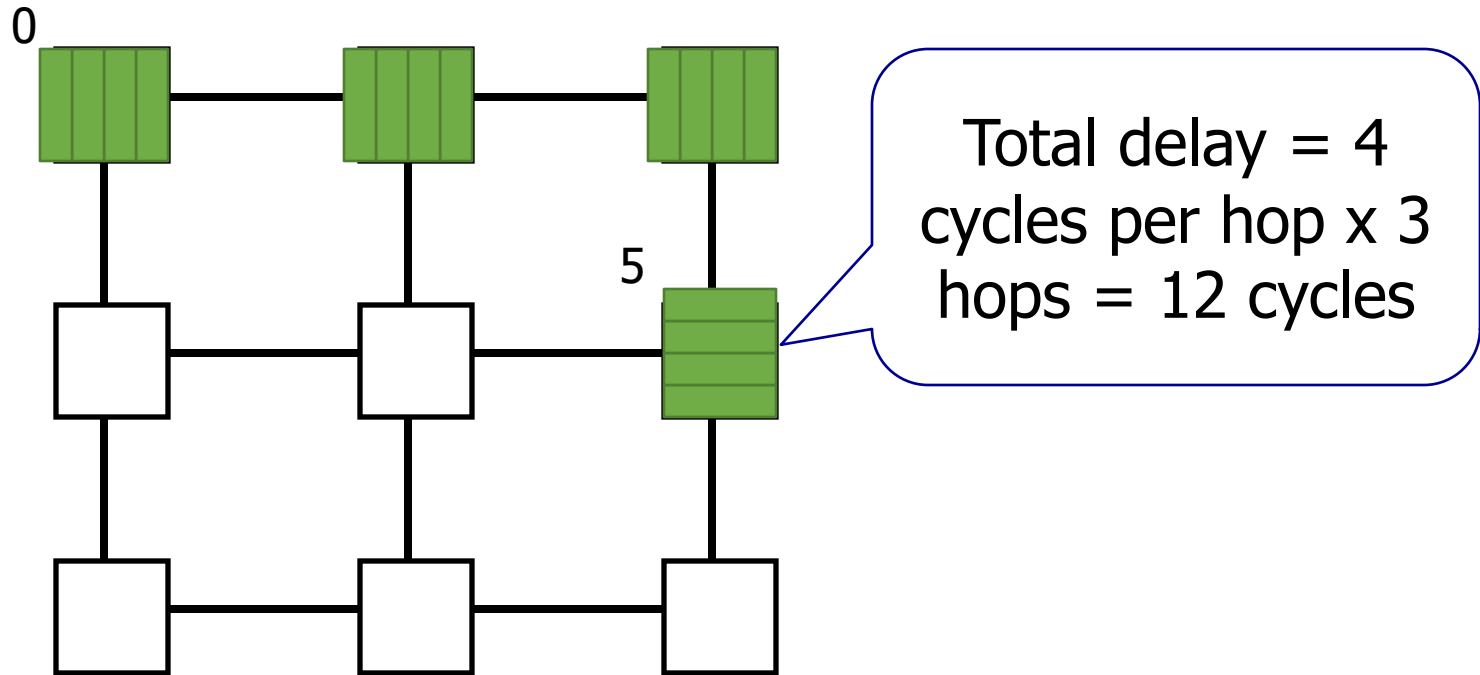
# PACKET-BASED FLOW CONTROL

- "Packet Switching"
  - Break messages into packets
  - Interleave packets on links
    - Better utilization
  - Requires per-node buffering to store packets in-flight waiting for output channel

- Two techniques
  - Store and Forward
  - Virtual Cut-Through

# PACKET-BASED: STORE AND FORWARD

- Links and buffers are allocated to **entire** packet

- Head flit **waits** at router until entire packet is received before being forwarded to the next hop
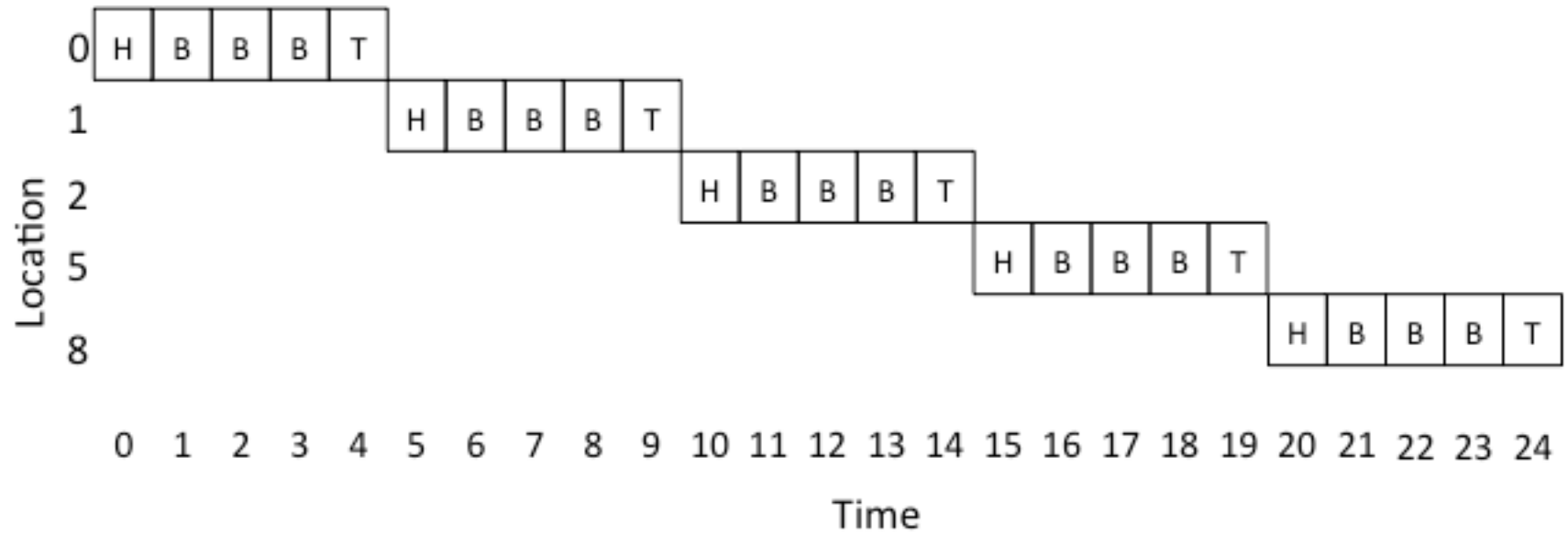
# STORE AND FORWARD EXAMPLE

0

5

Total delay = 4 cycles per hop x 3 hops = 12 cycles

**Not suitable on-chip.**
Why?

- High per-hop latency
  - Serialization delay paid at each hop
- Larger buffering required

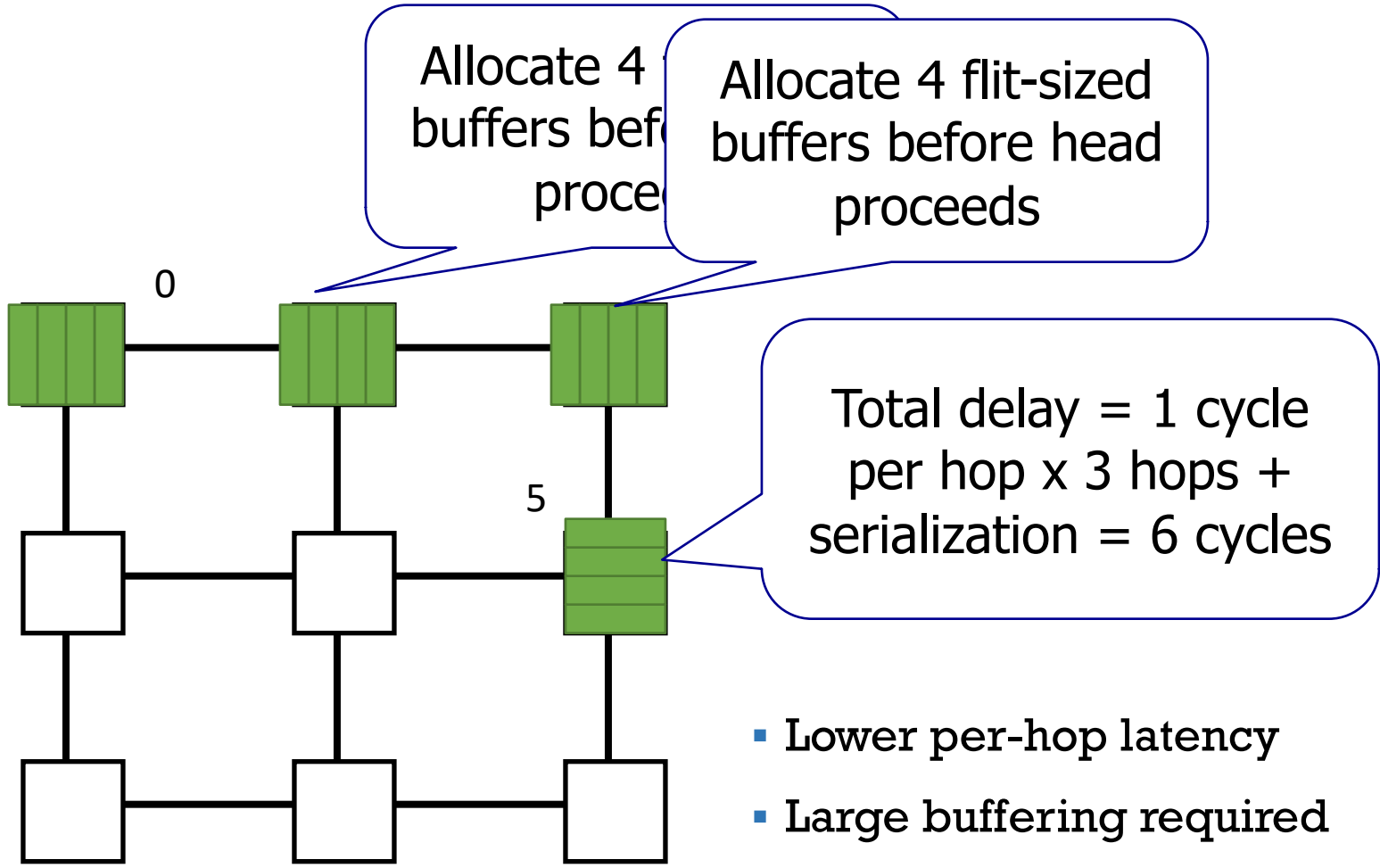# PACKET-BASED: VIRTUAL CUT-THROUGH

- Links and Buffers allocated to **entire** packets

- Flits can proceed to next hop before tail flit has been received by current router
  - But only if next router has enough buffer space for **entire** packet

# VIRTUAL CUT-THROUGH EXAMPLE

Allocate 4 flit-sized buffers before head proceeds

Allocate 4 flit-sized buffers before head proceeds

0

5

Total delay = 1 cycle per hop x 3 hops + serialization = 6 cycles

- Lower per-hop latency
- Large buffering required

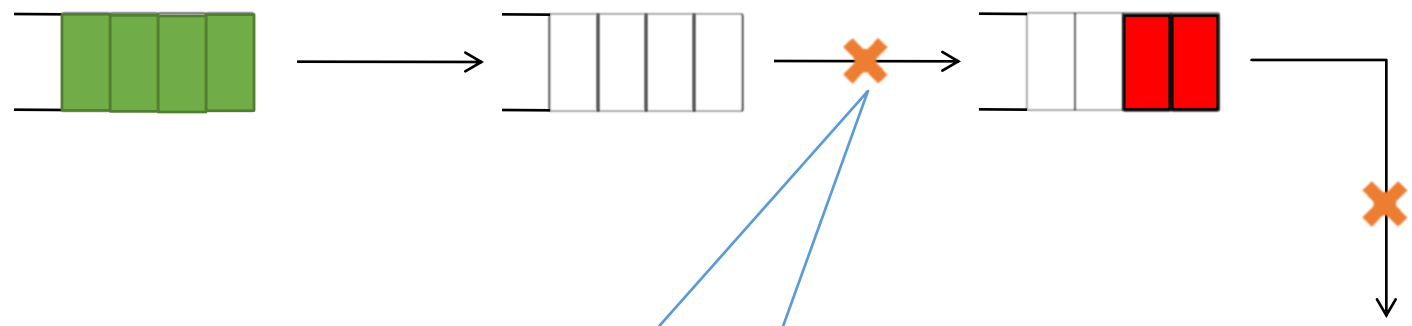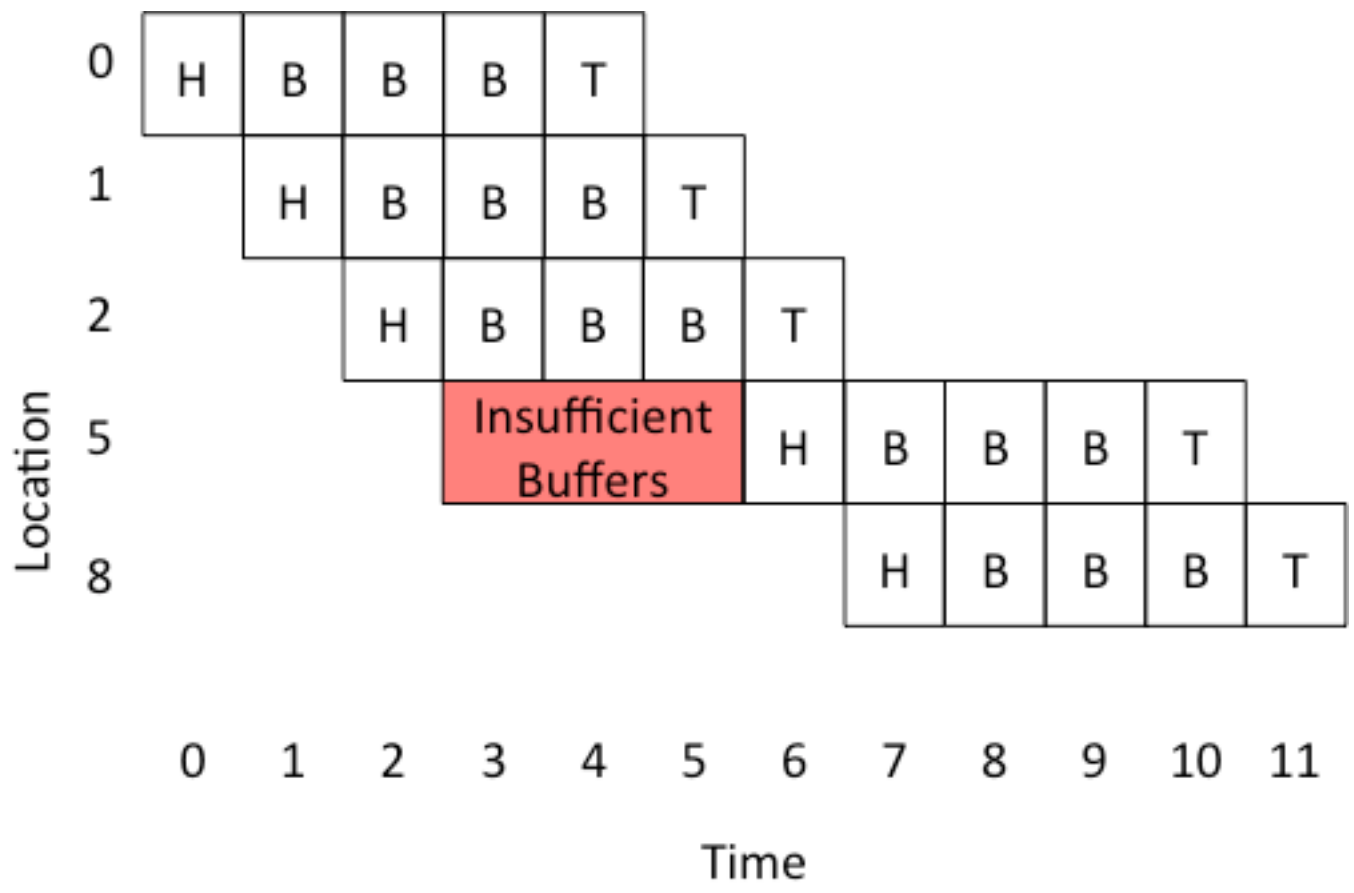# TIME-SPACE DIAGRAM: VIRTUAL CUT-THROUGH

# VIRTUAL CUT-THROUGH EXAMPLE (2)

Cannot proceed because only 2 flit buffers available

Throughput suffers from inefficient buffer allocation

# TIME-SPACE DIAGRAM: VIRTUAL CUT-THROUGH (2)

# FLIT-LEVEL FLOW CONTROL

- Like VCT, flit can proceed to next router before entire packet arrives
  - Unlike VCT, flit can proceed as soon as there is sufficient buffering for that **flit**

- Buffers allocated per flit rather than per packet
  - Routers do not need to have packet-sized buffers
  - Help routers meet tight area/power constraints

- Two techniques
  - Wormhole – link allocated per packet
  - Virtual Channel – link allocated per flit

# WORMHOLE FLOW CONTROL EXAMPLE

Red holds this channel: channel remains idle until red proceeds

Channel idle but red packet blocked behind blue

6 flit buffers/input port

Dest for Red

Buffer full: blue cannot proceed

Blocked by other packets

Dest for Blue

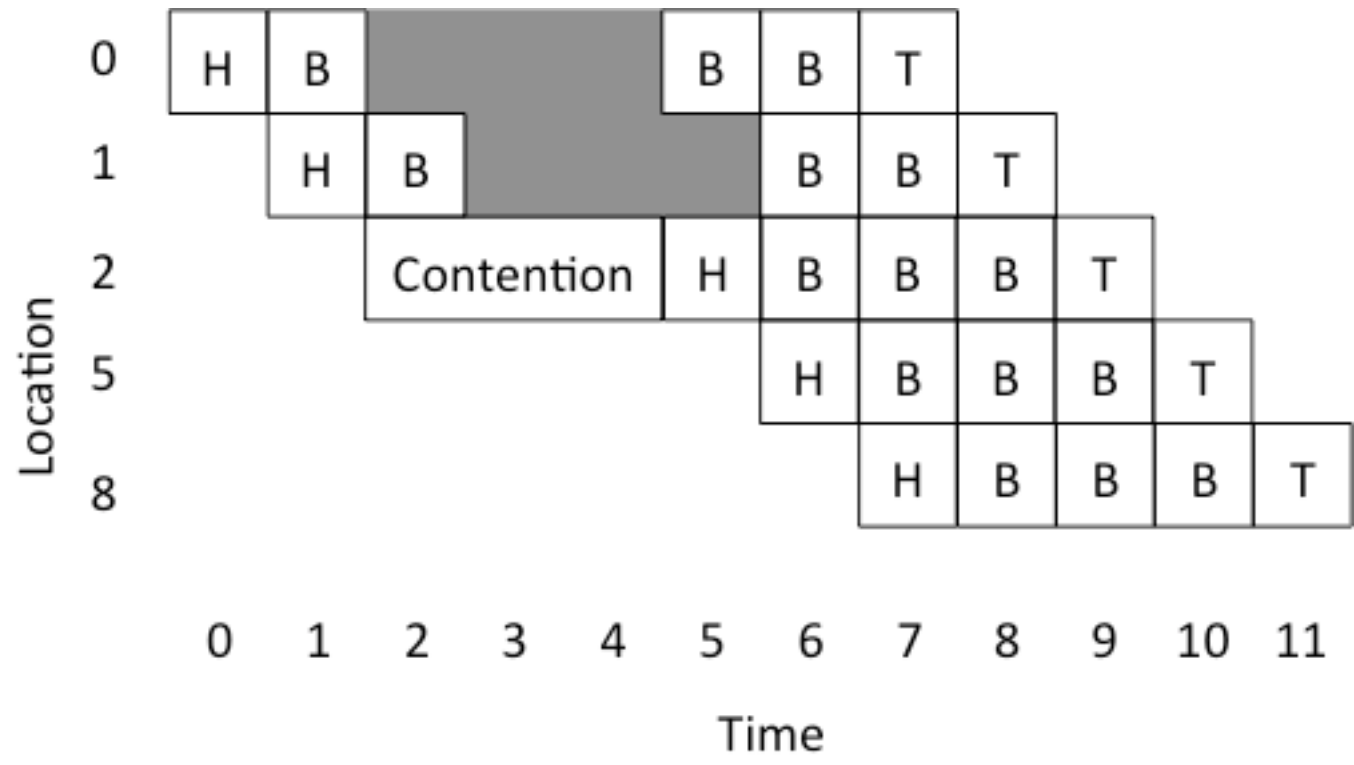## "Head-of-Line Blocking"

# WORMHOLE FLOW CONTROL

- Pros
  - More efficient buffer utilization (good for on-chip)
  - Low latency

- Cons
  - Poor link utilization: if head flit becomes blocked, all links spanning length of packet are idle
  - Cannot be re-allocated to different packet
  - Suffers from head of line (HOL) blocking

# TIME-SPACE DIAGRAM: WORMHOLE

# VIRTUAL CHANNEL FLOW CONTROL

- Like lanes on a highway
  - Flits on different VC can pass blocked packet
  - Link utilization improved

- Dual Use
  - Deadlock avoidance
  - Avoid Head-of-Line blocking

- Virtual channel implementation: multiple flit queues per input port
  - Share same physical link (channel)

# BLOCKING IN WORMHOLE FLOW CONTROL

virtual
channel

idle

chan p

idle

chan q

A

B

B

Node 1

Node 2

Node 3

Blocked

# VCS DECOUPLE DEPENDENCY BETWEEN BUFFER AND CHANNEL



Node 1    chan p    Node 2    chan q    Node 3

Blocked

# VIRTUAL CHANNEL FLOW CONTROL EXAMPLE



6 flit buffers/input port
3 flit buffers/VC

Dest for Red

Buffer full: blue cannot proceed

Blocked by other packets

Dest for Blue

## With Fair Interleaving

*Numbers under the buffers show number of flits in that VC's buffer, with capacity = 3.*

| In1 | AH | A1 | A2 | A3 | A4 | A5 |
|-----|----|----|----|----|----|----|
|     | 1  | 1  | 2  | 2  | 3  | 3  | 3 |

| In 2 | BH | B1 | B2 | B3 | B4 |  | B5 |
|------|----|----|----|----|----|--|----|
|      | 1  | 2  | 2  | 3  | 3  | 3 | 3 |

...

| Out |  | AH | BH | A1 | B1 | A2 | B2 |
|-----|--|----|----|----|----|----|----|

A Downstream     AH     A1     A2     A3

B Downstream        BH     B1     B2     B3

# TIME-SPACE DIAGRAM: VC FLOW CONTROL

## *With Fair Interleaving*

*Numbers under the buffers show number of flits in that VC's buffer, with capacity = 3.*



In1

| $A_H$ | $A_1$ | $A_2$ | $A_3$ | $A_4$ | $A_5$ | | $A_6$ | | $A_T$ |

1  1  2  2  3  3  3  3  3  3  3  2  2  1  1

In2

| $B_H$ | $B_1$ | $B_2$ | $B_3$ | $B_4$ | | $B_5$ | | $B_6$ | | $B_T$ |

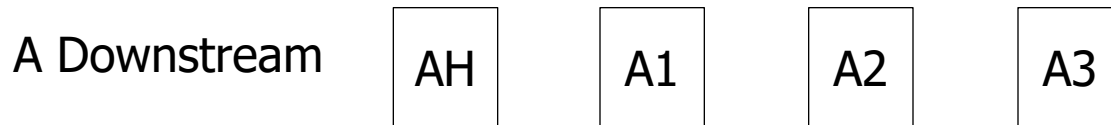1  2  2  3  3  3  3  3  3  3  3  3  2  2  1  1

Out

| $A_H$ | $B_H$ | $A_1$ | $B_1$ | $A_2$ | $B_2$ | $A_3$ | $B_3$ | $A_4$ | $B_4$ | $A_5$ | $B_5$ | $A_6$ | $B_6$ | $A_T$ | $B_T$ |

A downstream

| $A_H$ | $A_1$ | $A_2$ | $A_3$ | $A_4$ | $A_5$ | $A_6$ | $A_T$ |

B downstream

| $B_H$ | $B_1$ | $B_2$ | $B_3$ | $B_4$ | $B_5$ | $B_6$ | $B_T$ |

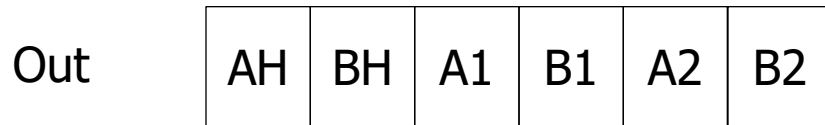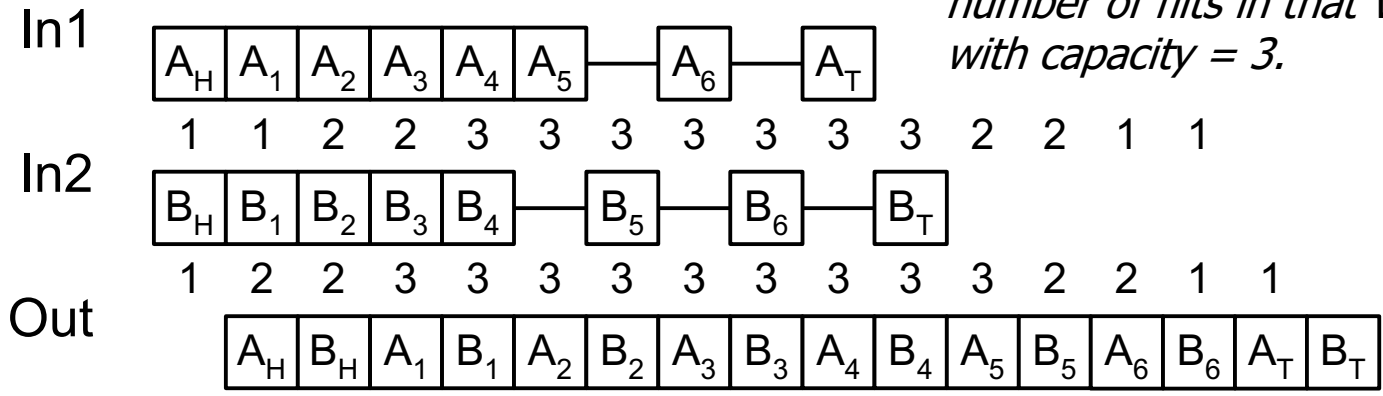**Latency of both packets got impeded due to fair interleaving!**

# TIME-SPACE DIAGRAM: VC FLOW CONTROL

## With Winner-Takes-All

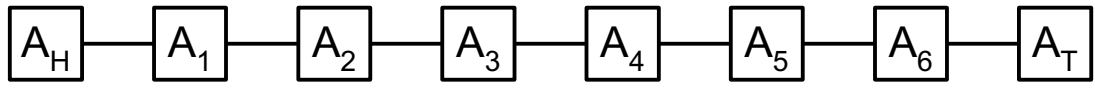*Numbers under the buffers show number of flits in that VC's buffer, with capacity = 3.*

In1

| $A_H$ | $A_1$ | $A_2$ | $A_3$ | $A_4$ | $A_5$ | $A_6$ | $A_T$ |
|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

In2

| $B_H$ | $B_1$ | $B_2$ | | | | | | | $B_3$ | $B_4$ | $B_5$ | $B_6$ | $B_T$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 2 | 1 |

Out

| $A_H$ | $A_1$ | $A_2$ | $A_3$ | $A_4$ | $A_5$ | $A_6$ | $A_T$ | $B_H$ | $B_1$ | $B_2$ | $B_3$ | $B_4$ | $B_5$ | $B_6$ | $B_T$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

A downstream

| $A_H$ | $A_1$ | $A_2$ | $A_3$ | $A_4$ | $A_5$ | $A_6$ | $A_T$ |
|---|---|---|---|---|---|---|---|

B downstream

| $B_H$ | $B_1$ | $B_2$ | $B_3$ | $B_4$ | $B_5$ | $B_6$ | $B_T$ |
|---|---|---|---|---|---|---|---|

Latency of packet A goes down by 7 cycles. (zero contention latency)
Latency of packet B is unaffected
(contention latency = serialization latency of packet A)

# SUMMARY OF TECHNIQUES

| | Links | Buffers | Comments |
|---|---|---|---|
| **Circuit-Switching** | Messages | N/A (buffer-less) | Setup & Ack |
| **Store and Forward** | Packet | Packet | Head flit waits for tail |
| **Virtual Cut Through** | Packet | Packet | Head can proceed |
| **Wormhole** | Packet | Flit | HOL |
| **Virtual Channel** | Flit | Flit | Interleave flits of different packets |

# DESIGNING A FLOW CONTROL PROTOCOL: MANAGING BUFFERS AND CONTENTION

# SUPPOSE WE HAVE A RING …



*For a Mesh, the analysis will be similar, with 5 ports (North, South, East, West, Core) instead of 2 (Ring, Core) ports*

# FLOW CONTROL PROTOCOL

**Input from Core**

?

**Ring**

**Output to Core**

**1.** Who should use output link?

**2.** What to do with the other flit (from ring/core)

*Have you seen this same situation in real life on a road network?*

# ROTARY



**1.** Who should use output link?

*Traffic already on ring has priority*

**2.** What to do with the other flit (from ring/core)

*Wait*

# FLOW CONTROL PROTOCOL

This is known as "arbitration"
The control structure is called an "arbiter"

Arbiter: Decides who uses the output link.

**Arbitration Result**
(Send input if no traffic on ring)

**Input from Core**          **Input from Core**



**Output to Core**          **Output to Core**

**Arbitration:**
Centralized or
Distributed?

**Centralized**
- Bus
- Crossbar

**Distributed**
- Ring
- Mesh

# FLOW CONTROL PROTOCOL

**3.** What should a flit do if its output is blocked?



**Input from Core**

New Flit

**Input from Core**

Full

**Output to Core**

**Output to Core**

# FLOW CONTROL OPTIONS

- What should a flit do if its output is blocked?
  - **Option 1:** Drop!
    - Send a NACK back for dropped packet or have a timeout
      - Source retransmits
      - Implicit congestion control
    - Flow control protocol on the Internet
    - **Advantage: can be bufferless!**
    - Challenges?
      - Latency and energy overhead of re-transmitting more than that of buffering so not preferred on-chip

# FLOW CONTROL OPTIONS

- **What should a flit do if its output is blocked?**
  - **Option 2:** Misroute!
    - As long as N input ports and N output ports, can send flit out of some other output port
      - called "bouncing" on a ring
    - **Advantage: can be bufferless!**
    - Challenges
      - Energy
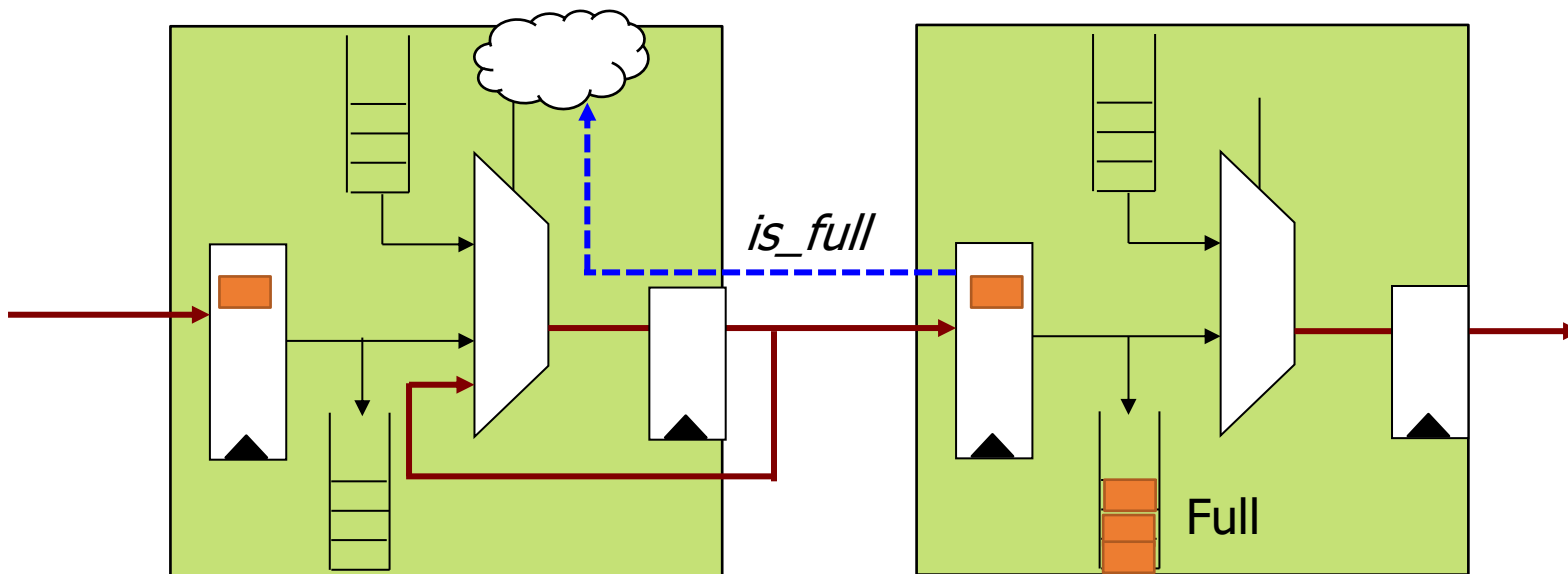        - Routes become non-minimal – more energy consumption at router latches and on links
      - Performance
        - Non-minimal routes – can lead to longer delays
      - Correctness
        - **Pt-to-Pt ordering violation** inside protocol
          - *Need mechanism to misroute subsequent packets from same source*
        - **Livelock!** – cannot *guarantee* forward progress
          - *Need to restrict number of misroutes of same packet*

# FLOW CONTROL OPTIONS

- What should a flit do if its output is blocked?
  - **Option 3:** Wait!
    - How? What about flit at previous router?
      - Signal back that it should wait too ("Backpressure")



*is_full*

Full

# ARBITRATION LOGIC

```
if (is_full)
    out = prev_out;
else if (in_ring.valid)
    out = in_ring;
else if (in_core.valid)
    out = in_core;
else
    out = 0;
```

*Note: if we use VC flow control, some other flit going into a VC that is not blocked can use the link*

# BACKPRESSURE SIGNALING MECHANISMS

- **On/Off Flow Control**
  - downstream router signals if it can receive or not

- **Credit-based Flow Control**
  - upstream router tracks the number of free buffers available at the downstream router

# ON/OFF FLOW CONTROL

- Downstream router sends a 1-bit on/off if it can receive or not
  - Upstream router sends only when it sees on

- Any potential challenge?
  - Delay of on/off signal
  - By the time the on/off signal reaches upstream, there might already be flits in flight
  - Need to send the off signal *once the number of buffers reaches a threshold* such that all potential in-flight flits have a free buffer

# ON/OFF TIMELINE WITH N BUFFERS

Node 1    Node 2

$N_{threshold}$ reached

t1

Flit

$N_{threshold}$ set to 3 to prevent flits departing Node 1 before t4 from overflowing

t2    Flit

Off    On/Off Delay = 2

t3    Flit

Process    Process Delay = 1

t4    Flit

Flit

Flit

Flit

$N_{total}$ set so that Node 2 does not run out of flits to send between t5 and t8

t5    Flit

On    $N_{threshold}$ +1 reached

t6    Flit

Process    Flit

t7    Flit

Flit    Flit

t8    Flit    Flit Delay = 1

Flit

# BACKPRESSURE SIGNALING MECHANISMS

- **On/Off Flow Control**
  - Pros
    - Low overhead: one-bit signal from downstream to upstream node, only switches when threshold crossed
  - Cons
    - Inefficient buffer utilization – have to design assuming worst case of $N_{threshold}$ flights in flight

# CREDIT-BASED FLOW CONTROL

- **Upstream router** tracks the **number of free buffers available at the downstream router**
  - Upstream router sends only if credits > 0

- When should credit be decremented at upstream router?
  - When a flit is sent to the downstream router

- When should credit be incremented at upstream router?
  - When a flit leaves the downstream router

# CREDIT TIMELINE

Node 0    Node 1    Node 2



t1  Credits=4    Credits=1    Credits=2

t2    Credit    Flit

Credits=0

t3  Process credit    Process Flit

t4  Credits=5    Credit    Flit

Credits=1

t5  Process credit    Credits=1

# BACKPRESSURE SIGNALING MECHANISMS

- **On/Off Flow Control**
  - Pros
    - Low overhead: one-bit signal
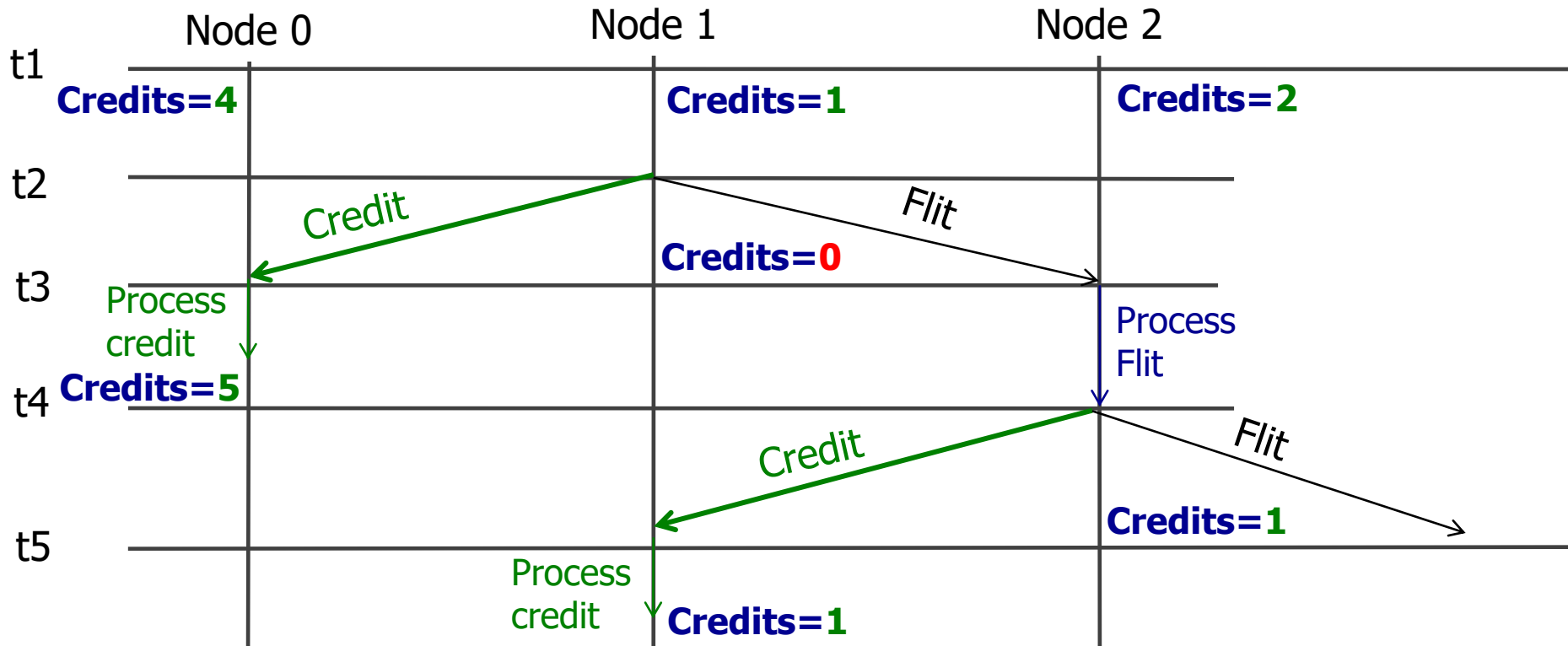  - Cons
    - Inefficient buffer utilization – have to design assuming worst case of $N_{threshold}$ flights in flight

- **Credit Flow Control**
  - Pros
    - Each buffer fully utilized - an keep sending till credits are zero (unlike on/off)
  - Cons
    - More signaling – need to signal upstream for every flit

# BACKPRESSURE AND BUFFER SIZING



**Node 0**     **Node 1**     **Node 2**

t1

**Credits=4**     **Credits=1**     **Credits=2**

t2

Credit     Flit

**Credits=0**

t3

Process credit    **Buffer Turnaround Time**    Process Flit

t4   **Credits=5**

or **Credit Round Trip Time**    Credit

**Credits=1**

t5

Process credit   **Credits=1**

*No flit can be sent into this buffer during this delay*

To prevent backpressure from limiting throughput, number of buffers >= turnaround time

# "BUFFER TURNAROUND TIME"

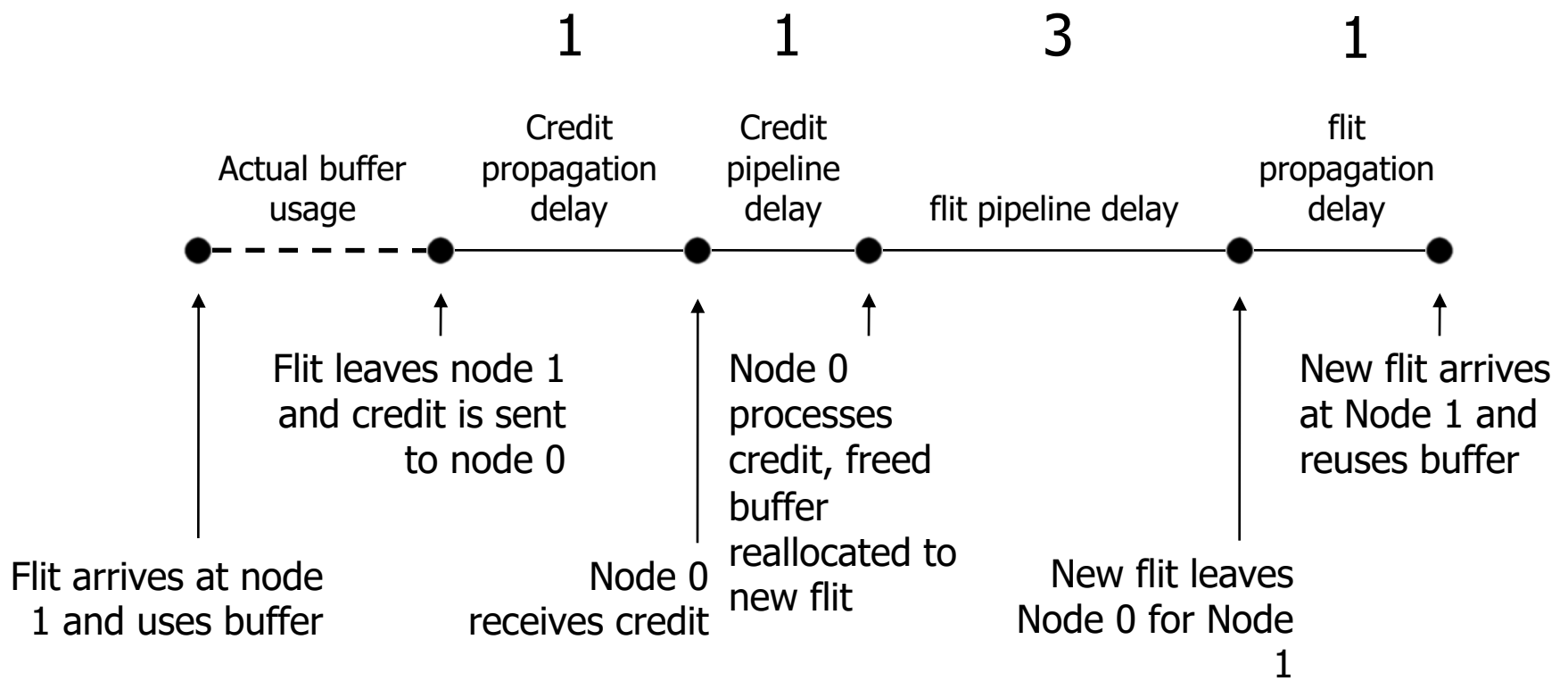| 1 | 1 | 3 | 1 |
|---|---|---|---|

**Actual buffer usage** — **Credit propagation delay** — **Credit pipeline delay** — **flit pipeline delay** — **flit propagation delay**

↑ Flit arrives at node 1 and uses buffer

↑ Flit leaves node 1 and credit is sent to node 0

↑ Node 0 receives credit

↑ Node 0 processes credit, freed buffer reallocated to new flit

↑ New flit leaves Node 0 for Node 1

↑ New flit arrives at Node 1 and reuses buffer
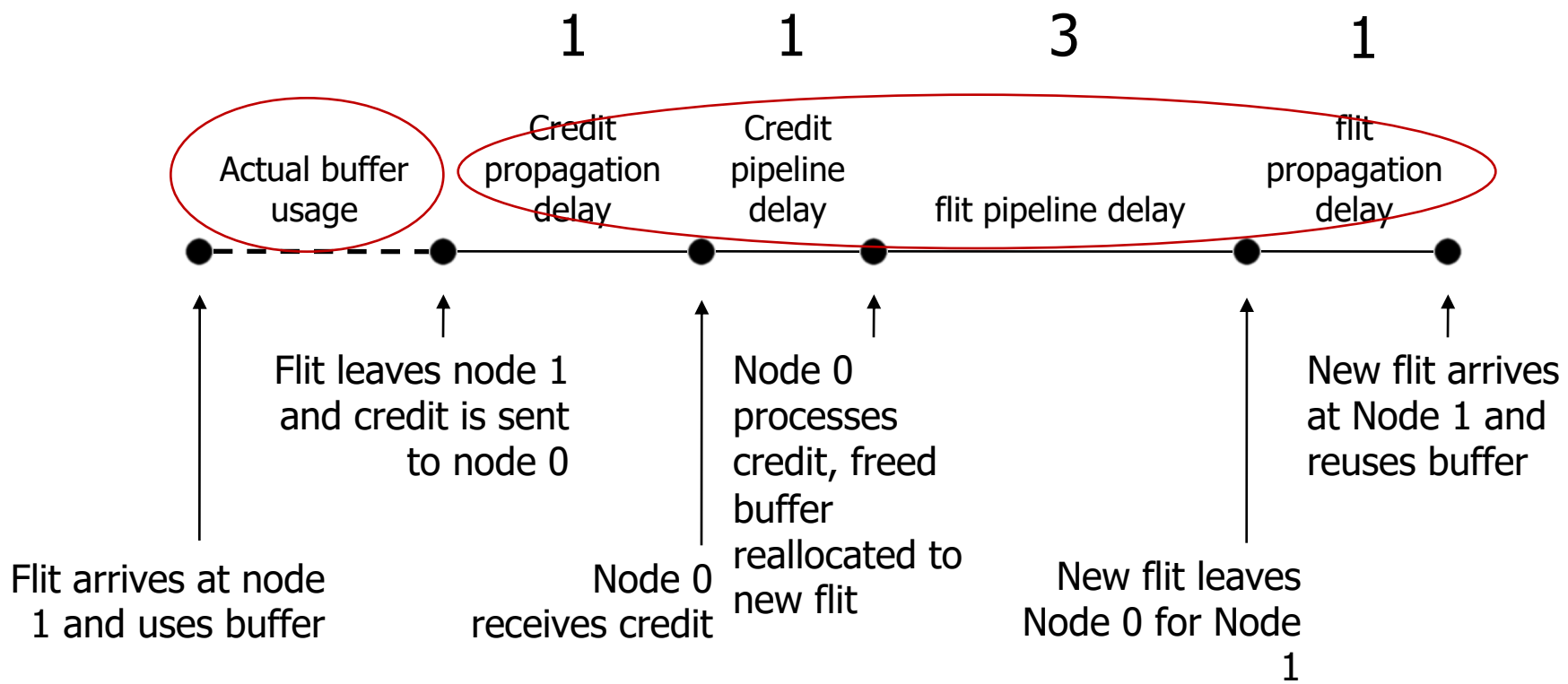
**How many buffers needed?**    1+1+3+1 = 6

**How many buffers needed in on/off flow-control?**    6 + 2 = 8
(off propogation + processing)

# BUT THIS IS INEFFICIENT

1          1          3          1

Actual buffer usage

Credit propagation delay

Credit pipeline delay

flit pipeline delay

flit propagation delay

Flit leaves node 1 and credit is sent to node 0

Node 0 processes credit, freed buffer reallocated to new flit

New flit arrives at Node 1 and reuses buffer

Flit arrives at node 1 and uses buffer

Node 0 receives credit

New flit leaves Node 0 for Node 1

*See: Flit Rsvn Flow Control, HPCA 2000*