

Lecture 15: NoC Innovations

- Today: power and performance innovations for NoCs

Network Power

- Power-Driven Design of Router Microarchitectures in On-Chip Networks, MICRO'03, Princeton

- Energy for a flit = $E_R \cdot H + E_{\text{wire}} \cdot D$
= $(E_{\text{buf}} + E_{\text{xbar}} + E_{\text{arb}}) \cdot H + E_{\text{wire}} \cdot D$

E_R = router energy

H = number of hops

E_{wire} = wire transmission energy

D = physical Manhattan distance

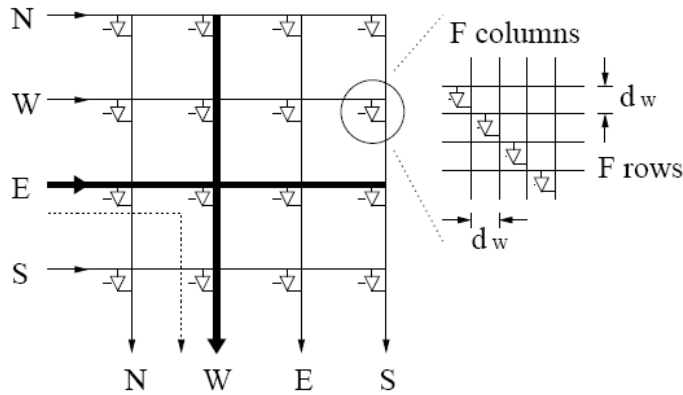
E_{buf} = router buffer energy

E_{xbar} = router crossbar energy

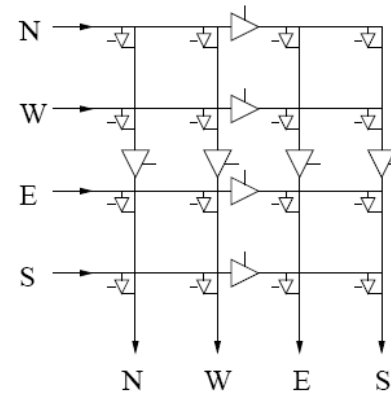
E_{arb} = router arbiter energy

- This paper assumes that $E_{\text{wire}} \cdot D$ is ideal network energy (assuming no change to the application and how it is mapped to physical nodes)

Segmented Crossbar



(a) A 4×4 matrix crossbar.

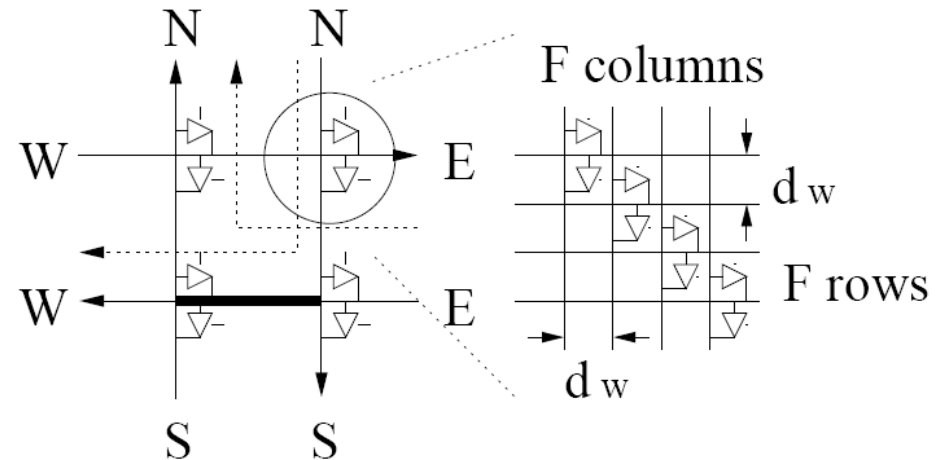


(b) A 4×4 segmented crossbar with 2 segments per line.

- By segmenting the row and column lines, parts of these lines need not switch \rightarrow less switching capacitance (especially if your output and input ports are close to the bottom-left in the figure above)
- Need a few additional control signals to activate the tri-state buffers
- Overall crossbar power savings: $\sim 15\text{-}30\%$

Cut-Through Crossbar

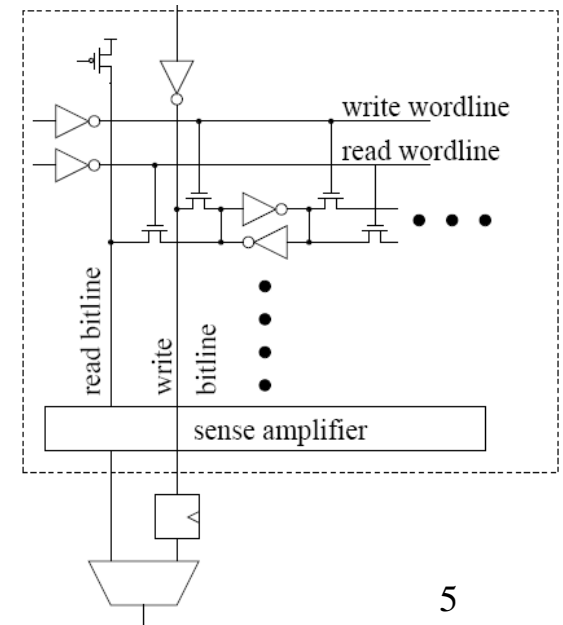
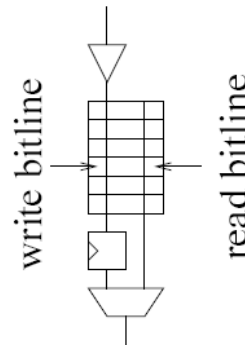
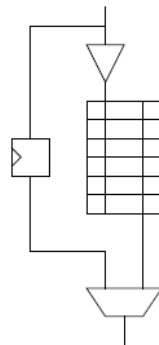
- Attempts to optimize the common case: in dimension-order routing, flits make up to one turn and usually travel straight
- $2/3^{\text{rd}}$ the number of tristate buffers and $1/2$ the number of data wires
- “Straight” traffic does not go thru tristate buffers
- Some combinations of turns are not allowed: such as $E \rightarrow N$ and $N \rightarrow W$ (note that such a combination cannot happen with dimension-order routing)
- Crossbar energy savings of 39-52%



(a) A 4×4 cut-through crossbar.

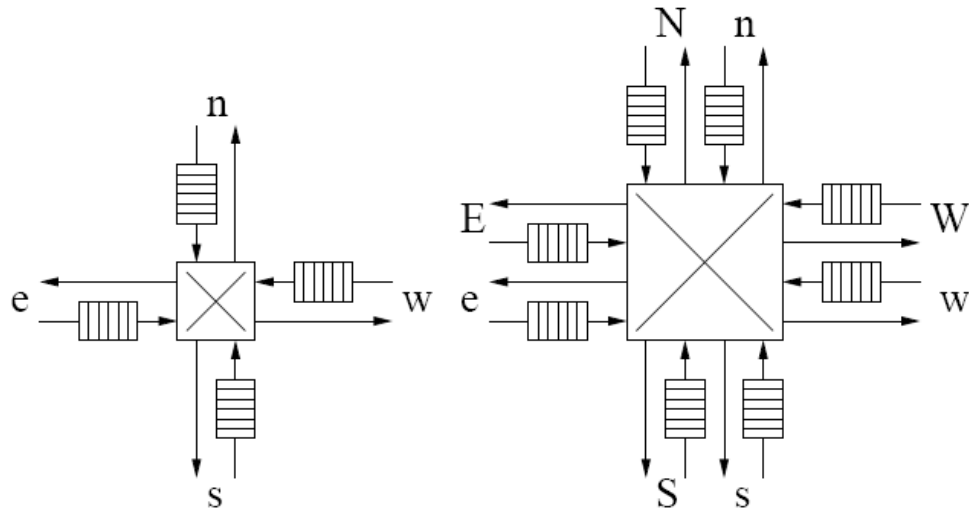
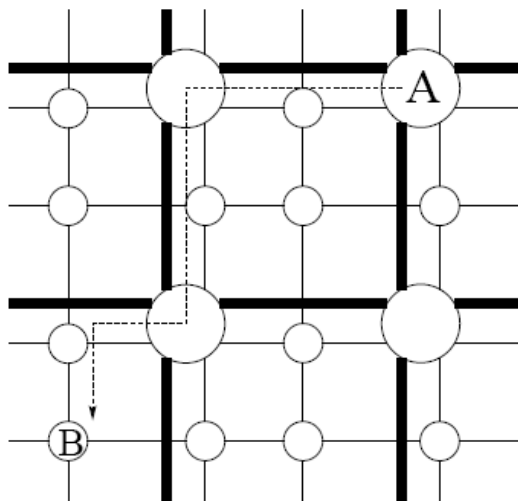
Write-Through Input Buffer

- Input flits must be buffered in case there is a conflict in a later pipeline stage
- If the queue is empty, the input flit can move straight to the next stage: helps avoid the buffer read
- To reduce the datapaths, the write bitlines can serve as the bypass path
- Power savings are a function of rd/wr energy ratios and probability of finding an empty queue



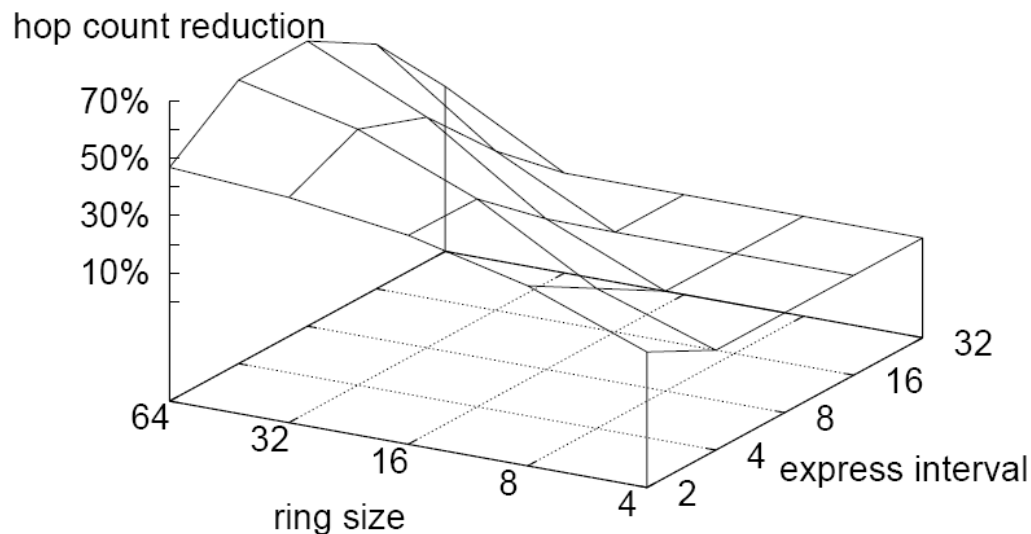
Express Channels

- Express channels connect non-adjacent nodes – flits traveling a long distance can use express channels for most of the way and navigate on local channels near the source/destination (like taking the freeway)
- Helps reduce the number of hops
- The router in each express node is much bigger now



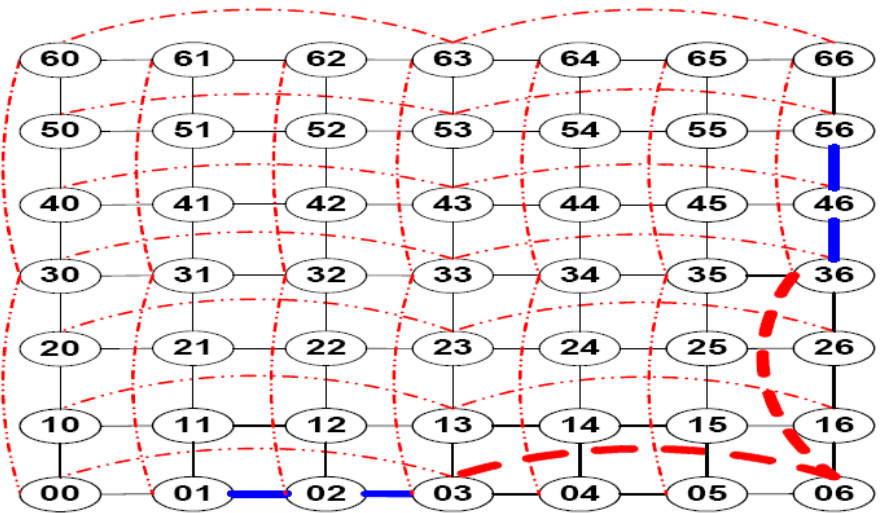
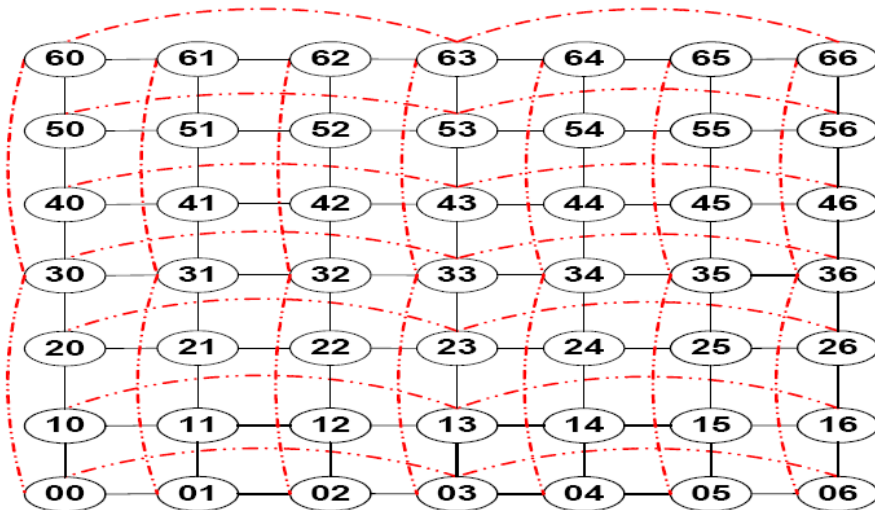
Express Channels

- Routing: in a ring, there are 5 possible routes and the best is chosen; in a torus, there are 17 possible routes
- A large express interval results in fewer savings because fewer messages exercise the express channels



Express Virtual Channels

- To a large extent, maintain the same physical structure as a conventional network (changes to be explained shortly)
- Some virtual channels are treated differently: they go through a different router pipeline and can effectively avoid most router overheads



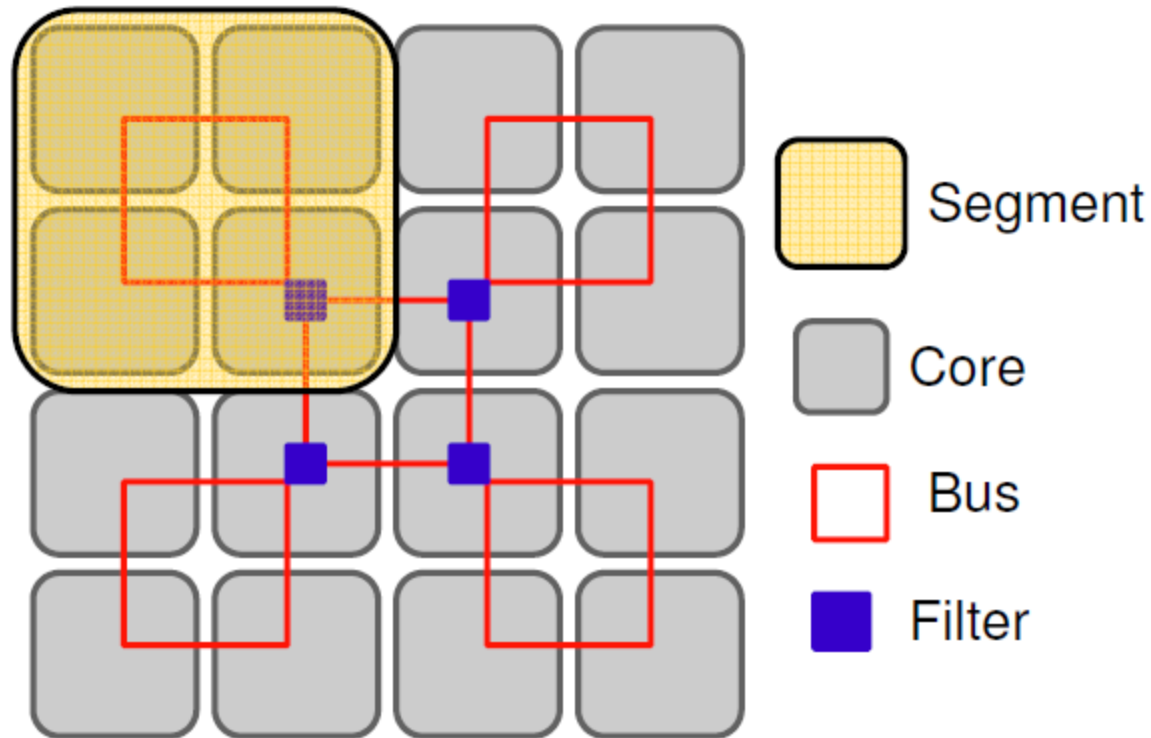
(b) VCs acquired from nodes 01 to 56

Router Pipelines

- If Normal VC (NVC):
 - at every router, must compete for the next VC and for the switch
 - will get buffered in case there is a conflict for VA/SA
- If EVC (at intermediate bypass router):
 - need not compete for VC (an EVC is a VC reserved across multiple routers)
 - similarly, the EVC is also guaranteed the switch (only 1 EVC can compete for an output physical channel)
 - since VA/SA are guaranteed to succeed, no need for buffering
 - simple router pipeline: incoming flit directly moves to ST stage
- If EVC (at EVC source/sink router):
 - must compete for VC/SA as in a conventional pipeline
 - before moving on, must confirm free buffer at next EVC router

Hierarchical Buses

Udipi et al., HPCA'10

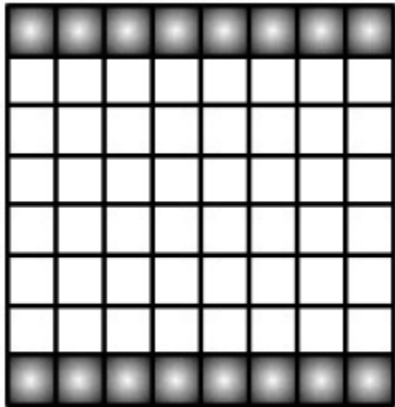


- Use buses to reduce routers
- Use bloom filters to stifle broadcasts
- Use page coloring to minimize travel beyond a bus

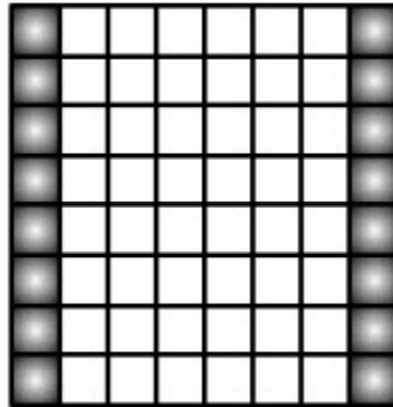
MC Placement

Abts et al., ISCA 2009

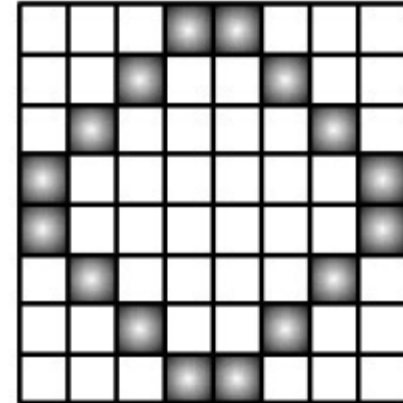
Tilera



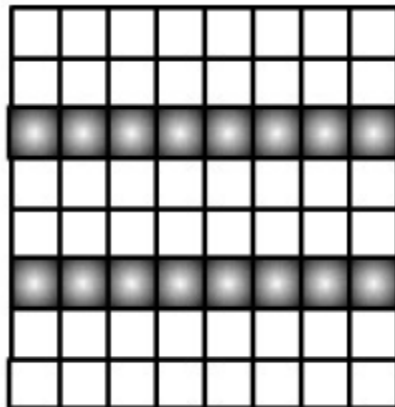
(a) row0_7



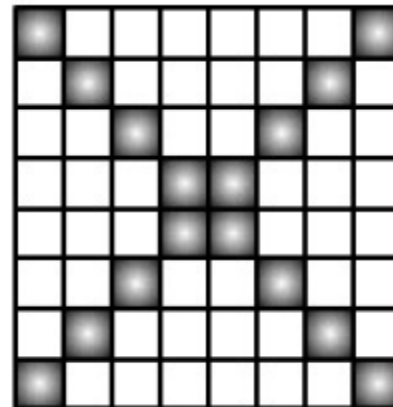
(b) col0_7



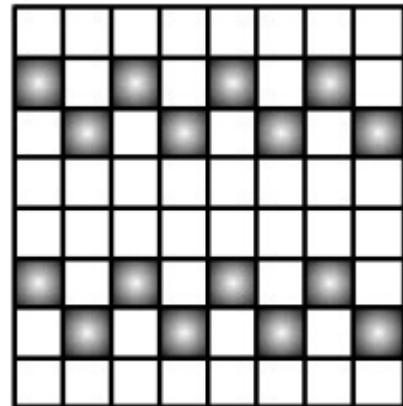
(e) diamond



(c) row2_5



(d) diagonal X



(f) checkerboard

MC Placement

- Diamond placement has the least link contention as it distributes traffic
- XY routing for requests and YX routing for replies is also effective in distributing traffic

- In speculative routers, two packets may be sent out at the same time on the link; receiver gets garbage
- Instead, send the XOR of the colliding packets
- In the next cycles, send the XOR of colliding packets, minus one of the colliding packets
- Receiver decodes all of the packets with XORs
- Saves one wasted cycle on every collision or mis-speculation

Title

- Bullet