

## HARAQ: Congestion-Aware Learning Model for Highly Adaptive Routing Algorithm in On-Chip Networks

Masoumeh Ebrahimi<sup>1</sup>, Masoud Daneshtalab<sup>1</sup>, Fahimeh Farahnakian<sup>1</sup>, Juha Plosila<sup>1</sup>, Pasi Liljeberg<sup>1</sup>, Maurizio Palesi<sup>2</sup>, Hannu Tenhunen<sup>1</sup>

<sup>1</sup>University of Turku, Finland, <sup>2</sup>University of Kore, Italy



## Outline

- Motivation
- The mad-y Method
- HARA: Highly Adaptive Routing Algorithm
- HARAQ: Q-Learning-based Highly Adaptive Routing Algorithm
- Results and Discussion
- Conclusion



## Motivation

- The occurrence of congestion in on-chip networks can severely degrade the performance due to increased message latency.
- Minimal Methods:
  - Minimal methods can propagate messages over two directions at each switch. When shortest paths are congested, sending more messages through them can deteriorate the congestion condition.
- Non-Minimal Methods:
  - Performance can severely deteriorate in non-minimal methods due to the uncertainty in finding an optimal path as they may choose longer paths and meanwhile delivering messages through congested regions.
- Reinforcement Learning / Q-routing Approach:
  - Due to maintaining large tables in each switch, the harware overhead is quite high.



## **Main Contributions**

- A low-restrictive non-minimal algorithm to provide several alternative paths between each pair of source and destination switches.
  - The algorithm uses only an extra virtual channel in the Y dimension.
  - Enables 180-degree turns on a single channel (i.e. a message can arrive through a channel that is previously used to deliver it).
- An efficient output selection strategy for finding a low-latency path from a source to a destination.
  - The output selection can efficiently estimate the latency of a message to reach its destination through each of the possible output channels.
  - Unlike typical Q-Routing methods, our proposed model is scalable and the size of Q-Tables is relatively small.



## The mad-y Method

- In 2D mesh-based network, three types of turns can be taken:
  - **0-degree**: a message transmits in a same direction with a possibility of switching between virtual channels. (0-degree-ch,0-degree-vc)
  - **90-degree:** a message transmits between the switches in perpendicular directions.
  - **180-degree turns (U turns):** a message is transferred to a channel in the opposite direction. (180-degree-vc,180-degree-ch)





## The mad-y Method

- To prove the deadlock freeness in mad-y, a two-digit number (a, b) is assigned to each output channel of a switch in n×m mesh network.
- 180-degree turns are not allowed in mad-y.





- As mad-y is a minimal adaptive routing method, it cannot fully utilize the eligible turns to route messages through less-congested areas.
- The aim of HARA, is to enhance the capability of the existing virtual channels in mad-y to reroute messages around congested areas and hotspots.
  - 180-degree turns are prohibited but can be incorporated in non-minimal routings.
    - One way to incorporate 180-degree turns is to examine the turns one by one to see whether the turn causes any cycle. After determining all allowable turns, in order to prove deadlock freeness, the numbering mechanism is utilized.
    - At first we use the numbering mechanism of the mad-y method to learn all 180-degree turns that can be taken in ascending order, and then modify the numbering mechanism to meet our requirements.







- In the non-minimal routing, employing only eligible turns at each switch is necessary but not sufficient to avoid blocking in the network.
  - The reason is that using the allowable turns, a message may not be able to find a path to the destination from the next hop and is blocked.
  - On the other hand, one of the aims of HARA is to fully utilize all eligible turns to present a low-restrictive adaptive method in the double-Y network.
- The output channels are selected in a way that not only the turn is allowable but also it is guaranteed that there is a path from the next switch to the destination.



## Potential output channels in HARA and Mad-y

Pos. InCh	Ν	S	Е	W	NE	NW	SE	SW
L	N1, N2	S2	Е	W	N1, N2, E	N1, W	S1, S2, E	S1, W
N1	-	S1, S2	Е	W	-	-	S1,S2, E	S1, W
N2	-	S2	Е	-	-	-	S2, E	-
<b>S1</b>	N1, N2	-	Е	W	N1, N2, E	N1, W	-	-
<b>S2</b>	N2	-	Е	-	N2, E	-	-	-
Ε	N1, N2	S2	-	W	N1, N2	N1, W	<u>S1, S2</u>	S1, W
W	N2	S2	Е	-	N2, E	-	S2, E	-

	Ν	S	Ε	W	NE	NW	SE	SW
L	N1, N2, S1,	N1, S1, S2,	N1, N2, S1,	N1, S1,	N1, N2, S1,	N1, S1,	N1, N2, S1,	N1, S1,
12	W	W	S2, E, W	W	S2, E, W	W	S2, E, W	W
N1	NO GI W	$\mathbf{S}1$ $\mathbf{S}2$ $\mathbf{W}$	N2, S1, S2,	S1 W	N2, S1, S2,	S1 W	N2, S1, S2,	S1 W
INI	1N2, 51, W	51, 52, W	E, W	51, W	E,W	51, W	E, W	51, W
N2	-	S2	S2, E	-	S2, E	-	S2, E	-
<b>C1</b>	N1, N2, S1,	N1, S1, S2,	N1, N2, S1,	N1, S1,	N1, N2, S1,	N1, S1,	N1, N2, S1,	N1, S1,
51	W	W	S2, E, W	W	S2, E, W	W	S2, E, W	W
<b>S2</b>	N2	-	N2, E	-	N2, E	-	N2, E	-
Б	N1, N2, S1,	N1, S1, S2,	N1, N2, S1,	N1, S1,	N1, N2, S1,	N1, S1,	N1, N2, S1,	N1,S1,
Ľ	W	W	S2, E, W	W	S2, E, W	W	S2, E, W	W
W	N2	S2	N2, S2, E	_	N2, S2, E	_	N2, S2, E	-



• HARA is deadlock-free and livelock-free.







## HARA → HARAQ: Q-Learning-based Selection





## MAIN CHARACTERISTICS OF HARAQ

#### **Q-Table Size:**

0

1

.

.

•

n

Num. of Switches





#### **Ouput Channels**

		N1	N2	S1	S2	Е	W
	Ν						
	S						
ns	Е						
tio	W						
osi	NE						
ď	NW						
	SE						
	SW						

Typical

C-routing

#### Proposed

Size\method	Q-Routing	C-Routing	R-Routing
8×8	128 bytes	40 bytes	24 bytes
16×16	512 bytes	64 bytes	24 bytes
32×32	2048 bytes	160 bytes	24 bytes



## MAIN CHARACTERISTICS OF HARAQ

### **Transferring Local and Global Information:**

- The congestion statuses are delivered over the channel whenever a message is transferred between two neighboring switches.
- A 4-bit congestion wire is used between each two neighboring switches to propagate local and global congestion information.
  - 2-bit local congestion information indicates the waiting time of a message from when the header flit is accommodated in an input buffer until an output channel is dedicated to it.
  - The global congestion information is a 4-bit value giving a global view of the latency from the output channel of the current switch to the destination switch region.
  - Upon connecting the input channel to the output channel, 2-bit local and 4-bit global values are aggregated into a 4-bit value and then transfer to the upstream switch.



## MAIN CHARACTERISTICS OF HARAQ

## **Table Initialization**

- Q-Routing models have an initial learning period during which it performs worse than minimal schemes.
  - The reason is that there is a possibility of choosing non-minimal paths even if the network is not congested.
  - All entries of Q-Tables are initialized such that minimal output channels are set to "0000" and non-minimal output channels are set to "1000" and never can be less than it.
  - In a low traffic condition, only minimal paths are selected while nonminimal paths are used to distribute traffic when the network gets congested.



## **Results and Discussion**

• We assess performance of HARAQ, a cycle-accurate NoC simulator developed in VHDL.

- DBAR and C-Routing schemes are implemented. For fairness, all methods utilize a fully adaptive routing function based on MAD-Y.
- The simulator inputs include the array size, the routing algorithm, the link width length, and the traffic type.
  - Wormhole switching
  - Data width is set to 64 bits
  - Frequncy 1GHz
  - Each input channel has a buffer (FIFO) size of 8 flits
  - The simulator is warmed up for 12,000 cycles and then the average performance is measured over another 200,000 cycles.
  - Two synthetic traffic profiles including uniform random and hotspot, along with SPLASH-2 [26] application traces are used.



Turun yliopisto University of Turku



Turun yliopisto University of Turku



Turun yliopisto University of Turku



## **Hardware Analysis**

- The whole platform of each scheme is synthesized by Synopsys Design Compiler.
- Each scheme includes switches, communication channels, & congestion wires.
- UMC 90nm technology at the operating frequency of 1GHz and supply voltage of 1V.
- We perform place-and-route, using Cadence Encounter, to have precise power and area estimations.
- The power dissipation of each scheme is calculated under the hotspot traffic profile near the saturation point (0.18) using Synopsys PrimePower in a 8×8 2D mesh

Network platforms	Area (mm <sup>2</sup> )	Avg. Power (W) dynamic & static	Max. Power (W) dynamic & static
DBAR	6.791	2.41	3.33
C-Routing	6.954	2.52	3.46
HARAQ	6.822	2.81	3.06



## Conclusion

- We proposed a highly adaptive routing algorithm based on minimal and non-minimal paths for on-chip networks.
- The presented algorithm provides a large number of paths for routing messages using only an extra virtual channel in the Y dimension.
- To choose a less congested path, we have utilized an optimized and scalable learning model to estimate the latency from each output channel to the destination switch.
- Results based on synthetic traffic profiles including uniform random and hotspot, along with SPLASH-2 application traces indicates that HARAQ performs better than DBAR and C-routing algorithms.



# Thank You