Odessa State Environmental University, Ukraine Sigma-T NVO, LTD, Dnipro, Ukraine

### Switching vs Routing within Multidimensional Torus Interconnect

Dmitry A. Zaitsev Serhii I. Tymchenko, Nataliia Z. Shtefan

http://daze.ho.ua

### **Torus interconnect**

- Three-dimensional torus network: IBM Blue Gene/L and Blue Gene/P, and the Cray XT3
- Five-dimensional torus network: IBM Blue Gene/Q
- Six-dimensional torus network: Fujitsu K computer, PRIMEHPC FX10 – threedimensional torus 3D mesh interconnect Tofu
- Fugaku, ~0.5 exaflops TOFU interconnect D

# Fugaku (Fujitsu, RIKEN)



http://top500.org





A64FX<sup>®</sup> Microarchitecture Manual <u>http://github.com/fujitsu/A64FX</u>

### A64FX block diagram



#### 6D Mesh/Torus Network

FUĴÎTSU

- Six coordinate axes: X, Y, Z, A, B, C
  - X, Y, Z: the size varies according to the system configuration
  - A, B, C: the size is fixed to 2×3×2
- Tofu stands for "torus fusion": (X, Y, Z) × (A, B, C)



### **Related Work**

- W.J. Dally, C.L. Seitz, "The torus routing chip," *Distrib. Comput.*, no. 1, 1986, pp. 187–196.
- A. Singh, W.J. Dally, B. Towles, and A.K. Gupta, "Globally Adaptive Load-Balanced Routing on Tori," *IEEE Computer Architecture Letters*, vol. 3, no. 1, pp. 2-2, Jan. 2004.
- P. Ren, M. A. Kinsy and N. Zheng, "Fault-Aware Load-Balancing Routing for 2D-Mesh and Torus On-Chip Network Topologies," *IEEE Transactions on Computers*, vol. 65, no. 3, pp. 873-887, 1 March 2016.
- P. Xie, H. Gu, K. Wang, X. Yu, S. Ma, "Mesh-of-Torus: a new topology for server-centric data center networks", *The Journal of Supercomputing*, vol. 75, pp. 255, 2019.

### **Torus Interconnect in Trend**

- Network-on-chip
- Interconnect of supercomputers and cluster
- Prospects for networks of service providers
- Prospects for campus and metropolitan networks
- Prospect for connecting Internet autonomous systems

### **Advantages of Torus Interconnect**

- Regular graph structure
- Short distance between nodes
- Many alternative (shortest) routes
- Possibility for packet delivery based on predefined switching rules
- Possibility of load balancing based on (random) alternative choice

### Neighborhoods in 2D torus





#### Von Neumann

Moore

### **Neighborhood of torus interconnect**

- Von Neumann neighborhood
- Mixed mesh and von Neumann interconnect
- Moore neighborhood is too dense in multidimensional space
- Generalized neighborhood is flexible, density is adjusted using a parameter
- Cross-By-Pass-Torus can by implemented as a generalized neighborhood with radius > 1

### **Basic Notions**

• Node address:

$$\mathfrak{i} = (i_0, i_1, \dots, i_{d-1}),$$

• Node ports:

$$p = (m, r)$$

• Neighbor node:

$$i_j \pm 1 \mod k$$

 Distance between nodes – Manhattan (taxicab) norm

$$D(i, i') = \sum_{j=0}^{d-1} \min(|i'_j - i_j|, k - |i'_j - i_j|)$$

# Rules based on the current and destination node addresses only



## Local switching rule

 Statement 1. Repeated forwarding a packet from a current node to the next node decreasing the distance (by unit), starting from the source node, provides finally the packet delivery to the destination node using one of the shortest paths.

### **Modifications of local switching rule**

- A) the first coordinate with nonzero difference (that corresponds to the deterministic dimension-order routing);
- B) random choice of coordinate among coordinates with nonzero difference;
- C) a random coordinate among coordinates with nonzero difference with the coordinate choice proportional to the coordinate differences.

# Rules taking into consideration the current node state



### **Modifications of local switching rule**

- D) the first coordinate with nonzero difference for *a free port*;
- E) random choice of coordinate among coordinates of *free ports* with nonzero difference;
- F) a random coordinate among coordinates of *free ports* with nonzero difference with the coordinate choice proportional to the coordinate differences.

### Example of Routes in (6,8)-torus

- Destination address: i' = (0,6,5,7,2,4)
- Current node address: i = (7,3,5,1,2,2)
- Port availability vector:  $a = \begin{pmatrix} 1 & 0 & 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 0 \end{pmatrix}$
- Address difference:  $\Delta = i' - i = (1,3,0,-2,0,2)$
- Distance: D(i, i') = 1 + 3 + 2 + 2 = 8

# Direction Preference for $\Delta = (1, 3, 0, -2, 0, 2)$

- 7  $\rightarrow$  0 before 7  $\rightarrow$  6  $\rightarrow$  5  $\rightarrow$  4  $\rightarrow$  3  $\rightarrow$  2  $\rightarrow$ 1  $\rightarrow$  0
- $3 \rightarrow 4 \rightarrow 5 \rightarrow 6$  before  $3 \rightarrow 2 \rightarrow 1 \rightarrow 0 \rightarrow 7 \rightarrow 6$
- 5
- $1 \rightarrow 0 \rightarrow 7$  before  $1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5 \rightarrow 6 \rightarrow 7$
- 2
- $2 \rightarrow 3 \rightarrow 4$  before  $2 \rightarrow 1 \rightarrow 0 \rightarrow 7 \rightarrow 6 \rightarrow 5 \rightarrow 4$

### **Port for Packet Forwarding**

- A) Port (0,1)
- B) Random choice between ports (0,1), (1,1), (3,-1), (5,1)
- C) Random choice between ports (0,1), (1,1), (3,-1), (5,1) with probabilities (<sup>1</sup>/<sub>8</sub>, <sup>3</sup>/<sub>8</sub>, <sup>1</sup>/<sub>4</sub>, <sup>1</sup>/<sub>4</sub>)
- D) Port (1,1)
- E) Random choice between ports (1,1),
   (3,-1)

# F)



# http://github/dazeorgacm/ts

```
daze@hare: ~/ts
<u>File Edit</u> View Search Terminal Help
daze@hare:~/ts$ ./ts --help
Simulator of traffic within d-dimensional torus of size k,
von Neuman neighborhood, local packet switching rules,
using shortest paths only (with random load balancing),
exponential distribution of time between packets,
node packet queue extraction: the first suitable.
USAGE: ts [options]
Options (keys):
 --d=dimension.
 --k=size.
 --r=rule: a-f.
 --cht=channel_time,
 --bl=buffer length,

    --lambda=node traffic intensity (exponential distribution).

 --maxst=halt simulation time,
 --dbg=debug_level, = 0,1,2...
Defaults: ts --d=3 --k=4 --r=a --cht=100 --bl=1000 --maxst=1000000 --dbg=0
```

### **Torus Simulator**

```
daze@hare: ~/ts
File Edit View Search Terminal Help
daze@hare:~/ts$ ./ts --d=6 --k=3 --lambda=0.01 --cht=200 --maxst=10000
***** Input information *****
torus dimensions d=6, size k=3
lambda=1.000000e-02, cht=200, bl=1000000
switching rule a
simulating...
***** Simulation Statistics *****
simulation time: 10001 (mtu)
generated packets: 73286
delevered packets: 64416
queued packets: 2984
dropped packets: 0 (0.000000e+00 %)
torus performanse: 6.440956e+00 (pkt/mtu)
torus load: 6.194488e+01 (%)
average hops per packet: 3.963503e+00
average packet c<u>h</u>annel time: 2.926989e+02 (mtu)
daze@hare:~/ts$
```

### **Create Packet Forwarding Canvas**



2D Torus 4x4 Example

### **Event Queue Records**

- i. Insertion of a packet into a node:
  - packet generation in the source node;
  - packet arrival to an intermediate or destination node.
- ii. Finishing of a packet transmission via a channel:
  - checking node queue for a packet waiting the corresponding port;
  - moving a packet from port to the next node inducing event (i) for the next node.

### **Basic Loop of Simulator**

- termination conditions checked and if one of them is true, the simulation process halts;
- the nearest time value among future events is chosen and the simulator current time is advanced to this moment of time;
- all the events having activation time equal to (or less than) the simulator current time are processed.

### **Evaluation of Actual Channel Transmission Time**



### Conclusions

- Torus structure: short distance between nodes, alternative shortest paths, possibility of packets delivery based on local rules
- 6 local rules offered including randomized choice for load balancing
- Torus Simulator *ts* has been implemented
- Simulation results show rather good performance and QoS

### **Basic References**

- D.A. Zaitsev, T.R. Shmeleva, and J.F. Groote, "Verification of Hypertorus Communication Grids by Infinite Petri Nets and Process Algebra," *IEEE/CAA Journal of Automatica Sinica*, 6(3), 2019, 733-742.
- T.R. Shmeleva, "Analysis of a Hypertorus Grid," In *Proc. of IEEE 38th International Conference Electronics and Nanotechnology ELNANO-2018*, Kyiv, Ukraine, April 24-26, 2018. NTUU "Igor Sikorsky Kyiv Polytechnic Institute", 2018. P. 56-59.
- D.A. Zaitsev, T.R. Shmeleva, W. Retschitzegger, B. Proll, "Security of grid structures under disguised traffic attacks," *Cluster Computing*, vol. 19, no. 3, 2016, pp. 1183-1200.
- D.A. Zaitsev, "Verification of Computing Grids with Special Edge Conditions by Infinite Petri Nets," *Automatic Control and Computer Sciences*, 2013, vol. 47, no. 7, pp. 403-412.
- D.A. Zaitsev, "A generalized neighborhood for cellular automata," *Theoretical Computer Science*, vol. 666, 2017, pp. 21-35,

