Shoal: A Network Architecture for Disaggregated Racks

Vishal Shrivastav (Cornell University)

Asaf Valadarsky (*Hebrew University of Jerusalem*) Hitesh Ballani, Paolo Costa (*Microsoft Research*) Ki Suh Lee (*Waltz Networks*) Han Wang (*Barefoot Networks*) Rachit Agarwal, Hakim Weatherspoon (*Cornell University*)

Traditional racks in datacenters



Disaggregated racks in datacenters



Disaggregated racks in datacenters



Challenges for disaggregated rack network

• Connect as many as an order of magnitude more nodes than traditional racks



Challenges for disaggregated rack network

• Connect as many as an order of magnitude more nodes than traditional racks



Challenges for disaggregated rack network

• Connect as many as an order of magnitude more nodes than traditional racks



Potential disaggregated rack network designs



Shoal is a network stack and fabric for disaggregated racks that is both low power and high performance (low latency, high throughput)

Key feature:

Shoal network fabric comprises purely *fast circuit switches* that can reconfigure within nanoseconds

Shoal is a network stack and fabric for disaggregated racks that is both low power and high performance (low latency, high throughput)

Key feature:

Shoal network fabric comprises purely *fast circuit switches* that can reconfigure within nanoseconds

Goal 1: Low power consumption



Packet switch

Ser Des Packet Processir g Ser Des Ser Des

Circuit switch

Circuit switches

- No buffering
- No packet processing
- □ No serialization/de-serialization

Consumes significantly less power than packet switches

Goal 2: High network performance

Key Challenge:

Need to explicitly set up circuits (reconfigure) before sending packets

Traditional circuit-switched networks

Uses switches with high reconfiguration delay, up to milliseconds

- Uses a central controller to decide the circuits (reconfiguration algorithm)
- □ Not suitable for low latency traffic

Shoal

Leverages circuit switches with nanosecond reconfiguration delay

Key Design Idea:

De-centralized, traffic agnostic reconfiguration algorithm

• Inspired from LB monolithic packet switches [Comp Comm'02]

Shoal for a single circuit switch network



Extending Shoal to a network of circuit switches

Time slot

	1	2	3	4	5	6	7
A	В	С	D	Е	F	G	Н
В	С	D	Ε	F	G	Η	А
С	D	Е	F	G	Н	А	В
D	Е	F	G	Н	Α	В	С
Е	F	G	Н	Α	В	С	D
F	G	Н	А	В	С	D	E
G	Н	А	В	С	D	Ε	F
Η	А	В	С	D	Ε	F	G



Extending Shoal to a network of circuit switches



	1	2	3	4	5	6	7
4	В	С	D	Е	F	G	Н
3	С	D	Е	F	G	Η	А
2	D	Е	F	G	Н	А	В
)	Е	F	G	Η	А	В	С
-	F	G	Η	А	В	С	D
:	G	Н	А	В	С	D	E
6	Н	А	В	С	D	Ε	F
1	Α	В	С	D	Ε	F	G



A non-blocking topology of circuit switches

Congestion in Shoal

Time slot 3 4 5 6 7 1 2 G Ε F Η Α В С D G Α F Н D Ε В С С G Н Α В Ε F D В С Η D F G Α Ε С G Н Α В D Ε F В С D Ε F G Н Α С G Ε F D Η В Α Η G Ε F D В С Α



Congestion control in Shoal



Each per-destination queue Q_i corresponding to destination i is bounded! $len(Q_i) \le 1 + incast_degree(i)$ packets

Key properties of Shoal

□ No central controller for reconfiguration

- □ Fully de-centralized, traffic agnostic reconfiguration logic
- □ Allows circuit switches to reconfigure at nanosecond timescales

Each per-destination queue in the network is bounded

Each packet traverses the network *at most* twice

- □ Worst-case 50% throughput compared to an ideal packet-switched network
- □ Can be compensated by allocating 2X bandwidth per node
- \Box Cost (Shoal) \leq Cost (packet-switched network with ½ bandwidth of Shoal)

Implementation

Stratix V FPGA Bluespec System Verilog

Implemented custom NIC and circuit switch on FPGA

Circuit switch implementation can reconfigure in < 6.4ns



Verified the queuing and throughput properties of Shoal on a 8-node testbed

Evaluation

Power consumption

For a 512-node rack

□ Packet-switched network comprises 24 64x50 Gbps packet switches

□ Shoal comprises 48 64x50 Gbps circuit switches

Packet-switched Network	8.72 KW	(58% of rack budget)
Shoal	2.55 KW	(17% of rack budget)

• Shoal consumes 3.5x less power than packet-switched network!

Evaluation

Network performance

- Packet-level simulator in C
- 512-node rack
- 5 disaggregated workload traces [OSDI'16]
- Shoal has 2X bandwidth (with comparable cost)
- Shoal performs comparable or better than several recent designs for packet-switched networks!



Short flows (0,100KB]

Long flows $[1MB,\infty)$

Conclusion



Thank you!

Shoal FPGA prototype and simulator code is available at: <u>https://github.com/vishal1303/Shoal</u>