



PURR: A Primitive for Reconfigurable Fast Reroute

Marco Chiesa KTH Royal Institute of Technology Code: <u>bitbucket.org/marchiesa/purr</u>





PURR: Thanks to my coauthors!

Andrzej, 16:10





Network resilience is a key yet challenging property



High resilience

Why is it that hard?

Routing reconvergence takes time!





Centralized controllers make reconvergence slower



Advanced SDN routing

e.g., intent-based logically centralized SDN



High resilience

failure occurs control-plane detection & route recomputation data-plane update

time



Hard to recover within 50ms!



Advanced SDN routing

e.g., intent-based logically centralized SDN



High resilience

"[carrier-grade networks] level of availability requires substantial **overprovisioning** and **fast reroute local recovery**, e.g., within 50 milliseconds" [1]

failure occurs
control-plane detection & route recomputation
data-plane update

[1] On low-latency-capable topologies, and their impact on the design of intra-domain routing. In SIGCOMM 2018

time



Fast Reroute (FRR): pre-computing failover paths





FRR entails solving two orthogonal problems:

1. control-plane: compute network-wide primary/backup forwarding rules



Fast Reroute (FRR): pre-computing failover paths



1. control-plane: compute **network-wide** primary/backup forwarding rules



FRR entails solving two orthogonal problems:

1. control-plane: compute network-wide primary/backup forwarding rules

2. data-plane: support conditional forwarding in each switch





PURR applies to P4 programmable switches





PURR applies to P4 programmable switches





First approach: packet recirculation



throughput reduction latency increase



FRR recirculation has high memory occupancy

Input				
FRR ₁ = 1 2 3 4				
FRR ₂ = 4 3 2 1				

match	action			
tag	fwd →	write & recirculate C		
1	1	tag := 2		
2	2	tag := 3		
3	3	tag := 4		
4	4	-		

throughput reduction latency increase high memory overhead port 4 fails port 3 fails



PURR: a primitive for reconfigurable FRR

A building block for implementing arbitrary FRR mechanisms



PURR: a primitive for reconfigurable FRR

A building block for implementing arbitrary FRR mechanisms

intriguing connection to (algorithmic) "string theory"









Input FRR₁ = 1 2 3 4





Input FRR₁ = 1 2 3 4

packet metadata	port status
FRR = 1	1111

match	action	
port status	fwd	
1 * * *	1	
* 1 * *	2	
* * 1 *	3	
* * * 1	4	



Input FRR₁ = 1 2 3 4



match	action	
port status	fwd	
1 * * *	1	
* 1 * *	2	
* * 1 *	3	
* * * 1	4	



Input FRR₁ = 1 2 3 4



match	action		
port status	fwd		
1 * * *	1		
* 1 * *	2		
* * 1 *	3		
* * * 1	4		



PURR: One tempting option: "Duplication" TCAM

Input FRR₁ = 1 2 3 4 FRR₂ = 2 3 4 1

packet metadata	port status				
FRR = 2	00 11				
similar to FRR recirculation:					
high memory overhead					

		match	action
	FRR	port status	fwd
	1	1 * * *	1
	1	* 1 * *	2
	1	* * 1 *	3
	1	* * * 1	4
	2	* 1 * *	2
	2	* * 1 *	3
	2	* * * 1	4
-	2	1 * * *	1

21



frr ports

 $b_1b_2b_3b_4b_5$

Input FRR₁ = 1 2 3 4 FRR₂ = 2 3 4 1

packet metadata	port status		
FRR = 1	1111		

Encoding FRR input:

- add a packet metadata field *frr_ports*
- map bits to the switch ports
- set bit to 1 to include a port
- set bit to 0 to skip a port



packet metadata port status			match		action	
FR	R = 1	1111		frr_ports	port status	fwd
			1 * * * *	1 * * *	1	
	match	action write	e	*1***	* 1 * *	2
>		frr_ports		* * 1 * *	* * 1 *	3
	FRR = 2	11110		* * * 1 *	* * * 1	4
	FRR = 2	01111		* * * * 1	1 * * *	1
bit-to-port mapping 12341 —						



packet metadata port status			match		action	
O FR	R = 1	1111		frr_ports	port status	fwd
			1 * * * *	1 * * *	1	
>	match	action write		*1***	* 1 * *	2
		frr_ports		* * 1 * *	* * 1 *	3
	FRR = 1	11110		* * * 1 *	* * * 1	4
	FRR = 2	01111		* * * * 1	1 * * *	1
bit-to-port mapping 12341 —						



packet metadata port status				ma	action			
FRR = 1		0111		frr_ports	port status	fwd		
				1 * * * *	1 * * *	1		
	match	action write		*1***	* 1 * *	2		
		frr_ports	Ν	* * 1 * *	* * 1 *	3		
	FRR = 1	11110		* * * 1 *	* * * 1	4		
	FRR = 2	01111		* * * * 1	1 * * *	1		
bit-to-port mapping 12341 —								



PURR: re-cycling TCAM entries

packet metadata port status				ma	action			
• FRR = 2		1111		frr_ports	port status	fwd		
				1 * * * *	1 * * *	1		
	match	action write		* 1 * * *	* 1 * *	2		
		frr_ports	>	* * 1 * *	* * 1 *	3		
	FRR = 1	11110		* * * 1 *	* * * 1	4		
	FRR = 2	01111		* * * * 1	1 * * *	1		
bit-to-port mapping 12341 —								



The key problem: How to compute the bit-to-port mapping that minimizes memory occupancy?

 $FRR_1 = 2310$ $FRR_2 = 0213$ $FRR_3 = 3021$ $FRR_4 = 1023$

Dit-to-port mapping =
$$2310213$$

FRR₁ = 2310
FRR₂ = 0213
FRR₃ = 3 021
FRR₄ = 102 3

Shortest Common Supersequence (SCS) problem without repetitions

- SCS without repetitions is **computationally hard** (based on [2])
- Dynamic Programming (**DPSCS**) computes optimum in exponential time



Fast-Greedy: a heuristic for solving this specific SCS

idea: remove the most frequent left-most element among the longest sequences

See the paper for:

 multi-table optimization







2020-10-04

Implementation feasibility

P4-based implementations:

 implemented different FRR mechanisms in P4 using PURR (e.g., F10 [1], arborescences [2], BFS, DFS, rotor router [3])



FPGA-based implementation:

implemented PURR on the NetFPGA-SUME platform

[1] V. Liu et al. "F10: A Fault-Tolerant Engineered Network" in NSDI 2013 [2] M. Chiesa et al. "On the Resiliency of Randomized Routing Against Multiple Edge Failures" in Transactions on Networking 2016 [3] Borokhovich et al ""Graph exploration algorithms" in HotSDN 2013







Evaluation: How does the FRR implementation impact memory and performance?

Two subquestions:

- 1. How much memory does PURR save?
- 2. How does performance in a datacenter vary depending on how one implements a FRR primitive?

See the paper for:

- multi-table optimization
- random vs tree-based FRR sequences
- FPGA chip occupancy
- low-size FRR sequences



2020-10-04

How much memory does PURR save? The "circular sequences" case

Input:

- switch with k ports
- 10 circular set of FRR sequences

"Duplication TCAM" FRR:

- k^2 number of TCAM entries
- 10 Top-of-Rack switches in a datacenter with F10 FRR [nsdi-13]
- 10 destinations with the "k arcdisjoint" FRR mechanism [ton-16]

With PURR encoding:

• *k-1* number of TCAM entries

[nsdi-13] V. Liu et al. "F10: A Fault-Tolerant Engineered Network" in NSDI 2013 [ton-16] M. Chiesa et al. "On the Resiliency of Randomized Routing Against Multiple Edge Failures" in Transactions on Networking 2016



2020-10-04

How much memory does PURR save? The "circular sequences" case

Input:

- switch with k ports
- 10 circular set of FRR sequences

For k = 24

- 92% less TCAM entries
- 470 instead of 5.760

"Duplication TCAM" FRR:

• k^2 number of TCAM entries

For k = 48

- 96% less TCAM entries
- 950 instead of 23.040

With PURR encoding:

• *k-1* number of TCAM entries

[nsdi-13] V. Liu et al. "F10: A Fault-Tolerant Engineered Network" in NSDI 2013 [ton-16] M. Chiesa et al. "On the Resiliency of Randomized Routing Against Multiple Edge Failures" in Transactions on Networking 2016



How much memory does PURR save? Fast-greedy performs close to the optimum



Input: randomly generated set of FRR sequences of length 7



How much memory Fast-greedy scales to

32 factorial possible FRR sequences The "duplication" TCAM or recirculation FRR approaches would not scale



Input: randomly generated set of FRR sequences



How does FCT vary depending on the FRR primitive? PURR improves both FCT and throughput



NS-3 simulations

Topology: 32-server Clos network, 10Gbps links

Workload: data-mining Transport: DCTCP One link failure at 0.5s



How does FCT vary depending on the FRR primitive? PURR improves both FCT and throughput



NS-3 simulations

Topology: 32-server Clos network, 10Gbps links

Workload: data-mining Transport: DCTCP One link failure at 0.5s



Conclusions: Keep calm and enjoy programmability

Fast Reroute is a critical functionality in today's network

• requires high throughput, low latency, fast reactiveness, small forwarding tables

P4 does not define an FRR built-in primitive

• pipeline compilers and control-plane must program the P4 pipeline



- PURR: We propose a lightweight TCAM-based FRR primitive
- an intriguing connection to algorithmic string theory
- no FRR-tailored hardware support
- improve performance by a factor of $\sim 2x$ w.r.t. FRR recirculation





Marco Chiesa KTH Royal Institute of Technology Code: <u>bitbucket.org/marchiesa/purr</u>



Backup slides





Smaller sequences on switches with high-density ports

Input:

- a 64-port programmable switch
- all possible three ports FRR sequences

FRR recirculation:

- #TCAM entries = 64*63*62 = 250K
- TCAM memory = 250K * (17+3) = 5Mb

"Duplication TCAM" approach:

- #TCAM entries = 64*63*62*3 = 750K
- TCAM memory = 750K * (20 + 3) = 17.2Mb

PURR encoding approach:

- #TCAM entries < (64*3) = 192
- TCAM memory < 192 * (64*4) = 50Kb

more than 99% memory reduction!



Three naïve approaches to implementing FRR

















match	action		match	action	>	match	action		match	action
FRR = 1	fwd(1)		FRR = 1	fwd(2)		FRR = 1	fwd(3)		FRR = 1	fwd(4)
FRR = 2	fwd(2)		FRR = 2	fwd(3)		FRR = 2	fwd(4)		FRR = 2	fwd(1)
FRR = 3	fwd(3)		FRR = 3	fwd(4)		FRR = 3	fwd(1)		FRR = 3	fwd(2)
FRR = 4	fwd(4)		FRR = 4	fwd(1)		FRR = 4	fwd(2)		FRR = 4	fwd(3)



A sequential search wastes hardware resources





PURR improves both FCT and throughput Two link failures, higher gains



NS-3 simulations

Topology: 32-server Clos network, 10Gbps links, $10\mu s$ link latency**Workload:** data-mining**Transport:** DCTCP



Evaluation: Random vs tree-based sequences





The web search workload



NS-3 simulations

Topology: 32-server Clos network, 10Gbps links, $10\mu s$ link latency **Workload:** web-search **Transport:** DCTCP