

Provable Data Plane Connectivity with Local Fast Failover

Introducing OpenFlow Graph Algorithms

Michael Borokhovich (Ben Gurion Uni, Israel)

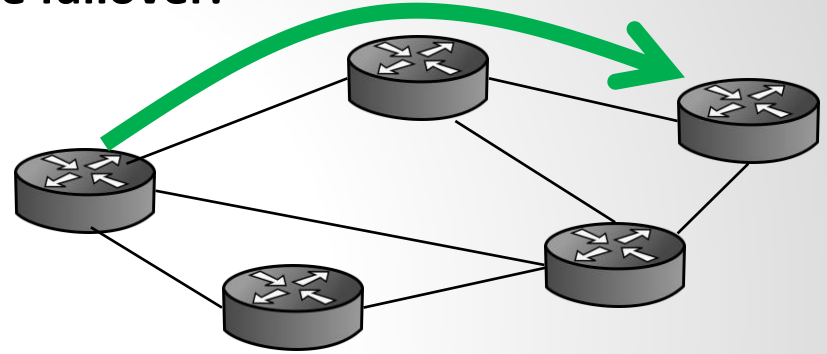
Liron Schiff (Tel Aviv Uni, Israel)

Stefan Schmid (TU Berlin & T-Labs, Germany)

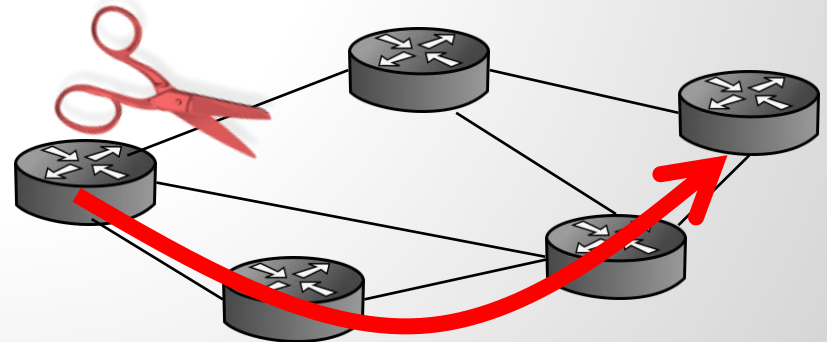
Robust Routing Mechanisms

- **Link failures** today are not uncommon [1]
- Modern networks provide **robust routing mechanisms**
 - i.e., routing which reacts to failures
 - example: MPLS local and global path protection

Before failover:

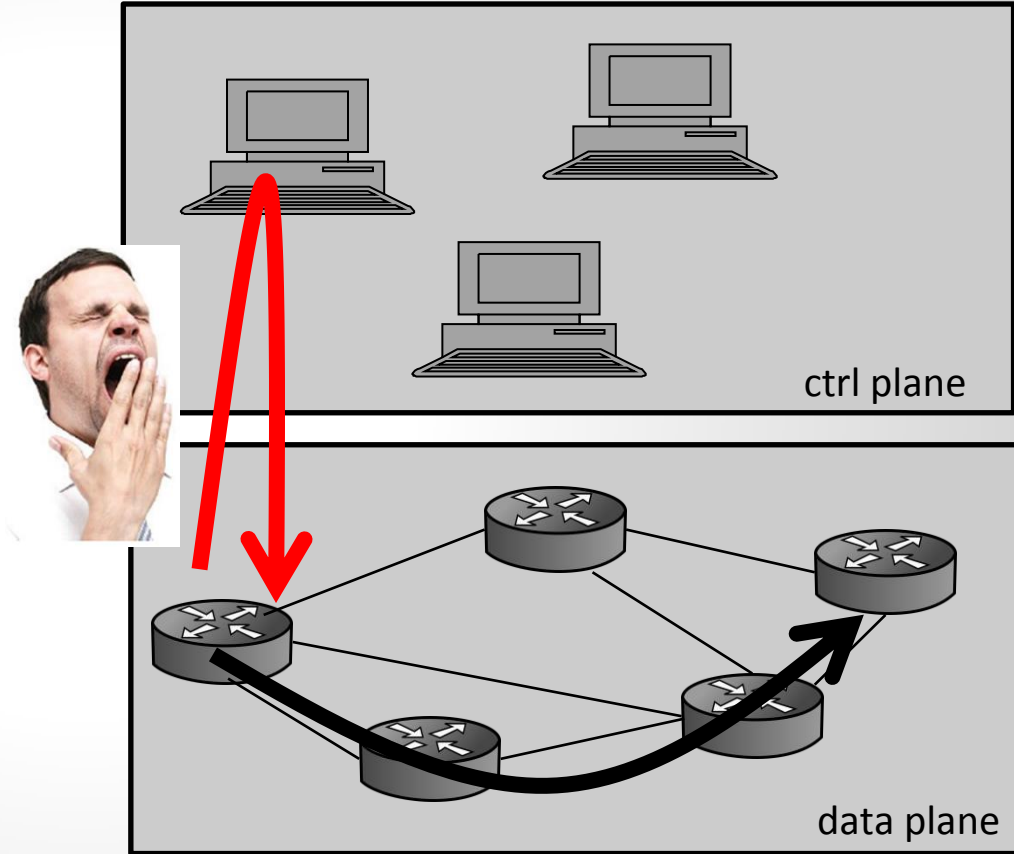


After failover:



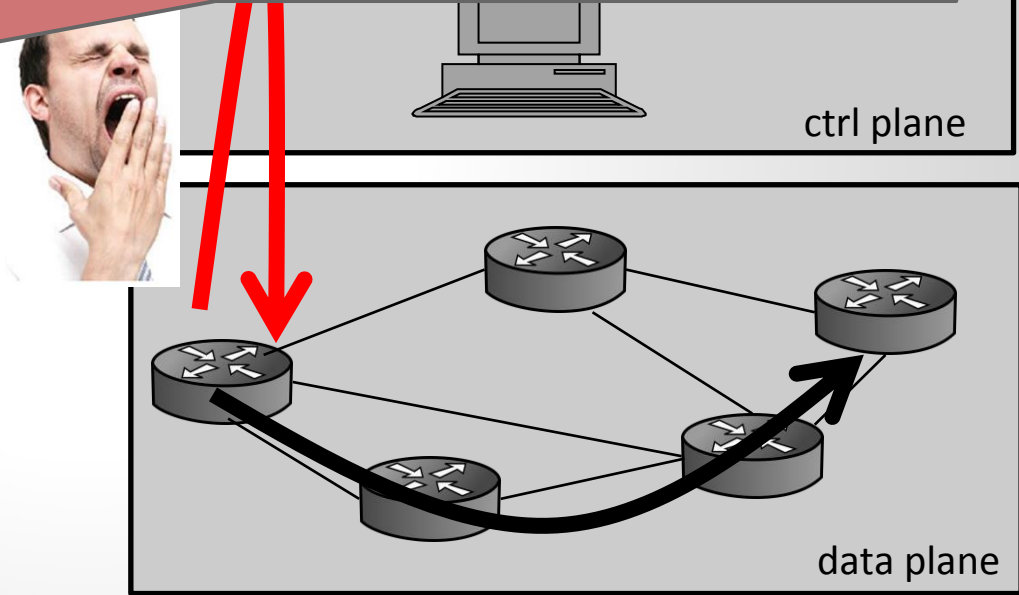
Fast In-band Failover

- Important that failover happens **fast = in-band**
 - Reaction time in control plane can be orders of magnitude slower [1]
- For this reason: **OpenFlow Local Fast Failover Mechanism**
 - Supports conditional forwarding rules (depend on the local state of the link: live or not?)
- Gives fast but local and perhaps **“suboptimal”** forwarding sets
 - Controller improves globally later...



However, not much is known about how to *use* the OpenFlow fast failover mechanism.
E.g.: **How many failures** can be tolerated without losing connectivity?

- Important that failover is **fast = in-band**
 - Reaction time in controller is orders of magnitude slower
- For this reason: **OpenFlow Local Fast Failover Mechanism**
 - Supports conditional forwarding rules (depend on the local state of the link: live or not?)
- Gives fast but local and perhaps **“suboptimal”** forwarding sets
 - Controller improves globally later...



However, not much is known about how to *use* the OpenFlow fast failover mechanism.
E.g.: **How many failures** can be tolerated without losing connectivity?

- Important that failover is **fast = in-band**

- Reaction time in control plane is orders of magnitude slower

- For this reason: **OpenFlow Local Fast Failover Mechanism**



ctrl plane

- Support for local fast failover (dependent on live network)
- Gives fast failover “suboptimal”
- Control plane

How to use mechanism is a **non-trivial problem** even if underlying network stays connected: (1) conditional failover rules need to be allocated **ahead of time**, without knowing actual failures, (2) views at runtime are **inherently local**.

How not to **shoot in your foot** with local fast failover (e.g., create forwarding loops)?

ne

Contribution: Very Robust Routing Possible with OpenFlow

Theorem: «Ideal» Forwarding Connectivity Possible

There exist algorithms which guarantee that packets always reach their destination, **independently of the number and locations** of failures, as long as the remaining network is connected.

Contribution: Very Robust Routing Possible with OpenFlow

Theorem: «Ideal» Forwarding Connectivity Possible

There exist algorithms which guarantee that packets always reach their destination, **independently of the number and locations** of failures, as long as the remaining network is connected.

Three algorithms:

- **Modulo**
- **Depth-First**
- **Breadth-First**

Essentially classic **graph algorithms** (routing, graph search) implemented **in OpenFlow**.

Make use of **tagging** to equip packets with meta-information to avoid forwarding loops.

Contribution: Very Robust Routing Possible with OpenFlow

Theorem: «Ideal» Forwarding Connectivity Possible

Analysis of their **complexity**: maximum stretch (route length compared to ideal route), number of tags, number of OpenFlow rules.

always reach
locations of
ed.

Three algorithms:

- **Modulo**
- **Depth-First**
- **Breadth-First**

Essentially classic **graph algorithms** (routing, graph search) implemented in **OpenFlow**.

Make use of **tagging** to equip packets with meta-information to avoid forwarding loops.

Overview of Contributions

High-Level Algorithms

Algorithm 2 Algorithm DFS

Input: current node: v_i , input port: in , packet dest: d ,
packet fallover global params: $pkt.start$, packet tag array:

Output: output
 $\{pkt.v_j\}_{j \in [n]}$

1: if $pkt.start =$

2: $out \leftarrow defc$

3: if out failed

4: $pkt.start$

5: $pkt.v_i, pa$

6: $out \leftarrow 1$

7: else

8: if $pkt.v_i, cu$

9: $pkt.v_i, pa$

10: $out \leftarrow pkt.s$

11: if $out = \Delta$

12: $out \leftarrow pk$

13: goto [19]

14: while out fai

15: $out \leftarrow out$

16: if $out = \Delta_i + 1$ when

17: $out \leftarrow pkt.v_i, par$

18: goto [19]

19: $pkt.v_i, cur \leftarrow out$

20: return out

Algorithm 1 Algorithm MOD

Input: current node: v_i , packet dest: d , packet tag array:

Output: output port: out

1: if no tag then {same as $pkt.v_i = 0$ }

2: $out \leftarrow default_route(i, d)$

3: else

4: $out \leftarrow (pkt.v_i \bmod \Delta_i) + 1$

5: $pkt.v_i \leftarrow out$

6: while out failed do

7: $out \leftarrow (pkt.v_i \bmod \Delta_i) + 1$

8: $pkt.v_i \leftarrow out$

9: return out

Algorithm 3 Algorithm BFS

Input: current node: v_i , input port: in , packet dest: d ,
packet fallover global params: $pkt.start$, packet tag array:

Output: output port: out

1: if $pkt.start = 0$ then

2: $out \leftarrow 1$

3: else

4: $out \leftarrow (pkt.v_i \bmod \Delta_i) + 1$

5: $pkt.v_i \leftarrow out$

6: while out failed do

7: $out \leftarrow (pkt.v_i \bmod \Delta_i) + 1$

8: $pkt.v_i \leftarrow out$

9: return out

Flow-Table Implementations

Flow Table Δ_{i-1}

Match	Instructions
0	Gr Δ_{i-1} , Table Δ_i
1	Drop

Flow Table Δ_i

Match	Instructions
0	Gr Δ_i , Table 0.2 (Send-Parent)
1	Drop

Table 5: DFS: Flow Tables of switch i .

Flow Table A (Start)

Match	Instructions	
$pkt.start$	dst	
0	1	Gr 0.1, Table 1
0	2	Gr 0.2, Table 1
...
0	n	Gr 0.n, Table 1
*	*	Table B

Flow Table B

Match			Instructions
in	$pkt.v_i, cur$	$pkt.v_i, par$	
1	0	*	$pkt.v_i, par \leftarrow in$, Table 2
1	1	2	Table 3
2	2	3	Table 4
3	3	4	Table 5
...
$\Delta_i - 2$	$\Delta_i - 2$	$\Delta_i - 1$	Table Δ_i
$\Delta_i - 1$	$\Delta_i - 1$	Δ_i	Table C
*	0	*	$pkt.v_i, par \leftarrow in$, Table 1
1	1	*	Table 2
2	2	*	Table 3
3	3	*	Table 4
...
$\Delta_i - 1$	$\Delta_i - 1$	*	Table Δ_i
Δ_i	Δ_i	*	Table C
1	*	*	Fwd 1
2	*	*	Fwd 2

Complexity Analysis

THEOREM 1. MOD ensures data plane connectivity when-

Algorithm	Packet Memory	Message count	Rules space
Module	$n \log d$	$\text{Exp}(n)$	$O(n^d)$
DFS	$n \log d$	$2 E $	$O(n^d)$
BFS	$n \log d$	$2kn$	$O(n^d)$
BFS*	$k(\log d + \log n)$	$2kn$	$O(n^{(d+k)})$

Related Work

- Borokhovich, OPODIS'13
- [1] Liu et al. NSDI'13
- Graph-search literature

Overview of Contributions

We expect that our algorithms scale up to **500-node networks** (ignoring link capacities) (e.g., using our NoviKit 250 switches, with 32MB flow table space and full support for extended match fields)

Flow-Table Implementations

```

2: out ← out + 1
3: if out failed do
4:   pkt.start ← pkt.start + 1
5:   pkt.vi, par ← pkt.vi, par + 1
6:   out ← 1
7: else
8:   if pkt.vi, cur mod Δi + 1
9:     pkt.vi, cur ← pkt.vi, cur + 1
10:  out ← pkt.vi, cur
11:  if out = 0 then
12:    out ← pkt.vi, cur mod Δi + 1
13:  got ← out
14:  while out = 0 do
15:    out ← out + 1
16:  if out = 0 then
17:    out ← out + 1
18:  if out = 0 then
19:    pkt.start ← pkt.start + 1
20:  return out
27: if out = 0 then
28:   out ← 1
29:  while out failed do
30:    out ← out + 1
31:  if out = Δi + 1 then
32:    Drop
33:  pkt.vi, cur ← out
34:  return out
    
```

Flow Table Δ_i

Match	Instructions
0	Gr Δ_i , Table 0.2 (Send-Parent)
1	Drop

Table 5: DFS: Flow Tables of switch i .

Flow Table A (Start)

Match	Instructions
0	Gr 0.1, Table 1
0	Gr 0.2, Table 1
...	...
0	Gr 0.n, Table 1
*	* Table B

Flow Table B

Match			Instructions
in	$pkt.vi, cur$	$pkt.vi, par$	
1	0	*	$pkt.vi, par \leftarrow in$, Table 2
1	1	2	Table 3
2	2	3	Table 4
3	3	4	Table 5
...
$\Delta_i - 2$	$\Delta_i - 2$	$\Delta_i - 1$	Table Δ_i
$\Delta_i - 1$	$\Delta_i - 1$	Δ_i	Table C
*	0	*	$pkt.vi, par \leftarrow in$, Table 1
1	1	*	Table 2
2	2	*	Table 3
3	3	*	Table 4
...
$\Delta_i - 1$	$\Delta_i - 1$	*	Table Δ_i
1	Δ_i	*	Table C
1	*	*	Fwd 1
2	*	*	Fwd 2

Complexity Analysis

THEOREM 1. MOD ensures data plane connectivity when-

Algorithm	Packet Memory	Message count	Rules space
Module	$n \log d$	$\text{Exp}(n)$	$O(n^d)$
DFS	$n \log d$	$2 E $	$O(n^d)$
BFS	$n \log d$	$2kn$	$O(n^d)$
BFS*	$k(\log d + \log n)$	$2kn$	$O(n^{(d+k)})$

Related Work

- Borokhovich, OPODIS'13
- [1] Liu et al. NSDI'13
- Graph-search literature

Overview of Contributions

We expect that our algorithms scale up to **500-node networks** (ignoring link capacities) (e.g., using our NoviKit 250 switches, with 32MB flow table space and full support for extended match fields)

Same objective: ideal connectivity. But their **link-reversal algorithms** not applicable to OpenFlow: require dynamic state at router.

Inherent tradeoffs between **robustness and network load** of failover **without tagging**.

Complexity Analysis

Lower bounds with **implications on optimality** of our algorithms.

Algorithm	Rules space
Module	$O(n^d)$
DFS	$O(n^d)$
BFS	$O(n^d)$
BFS*	$O(n^{(d+k)})$

Related Work

- Borokhovich, OPODIS'13
- [1] Liu et al. NSDI'13
- Graph-search literature

Conclusion

- Fast failover: example of a function that should be kept **in the data plane**



- Our result shows that **non-trivial functions** can be computed in the OpenFlow data plane!
- Our algorithms: may serve in compilers for **higher-level languages**, e.g., FatTire