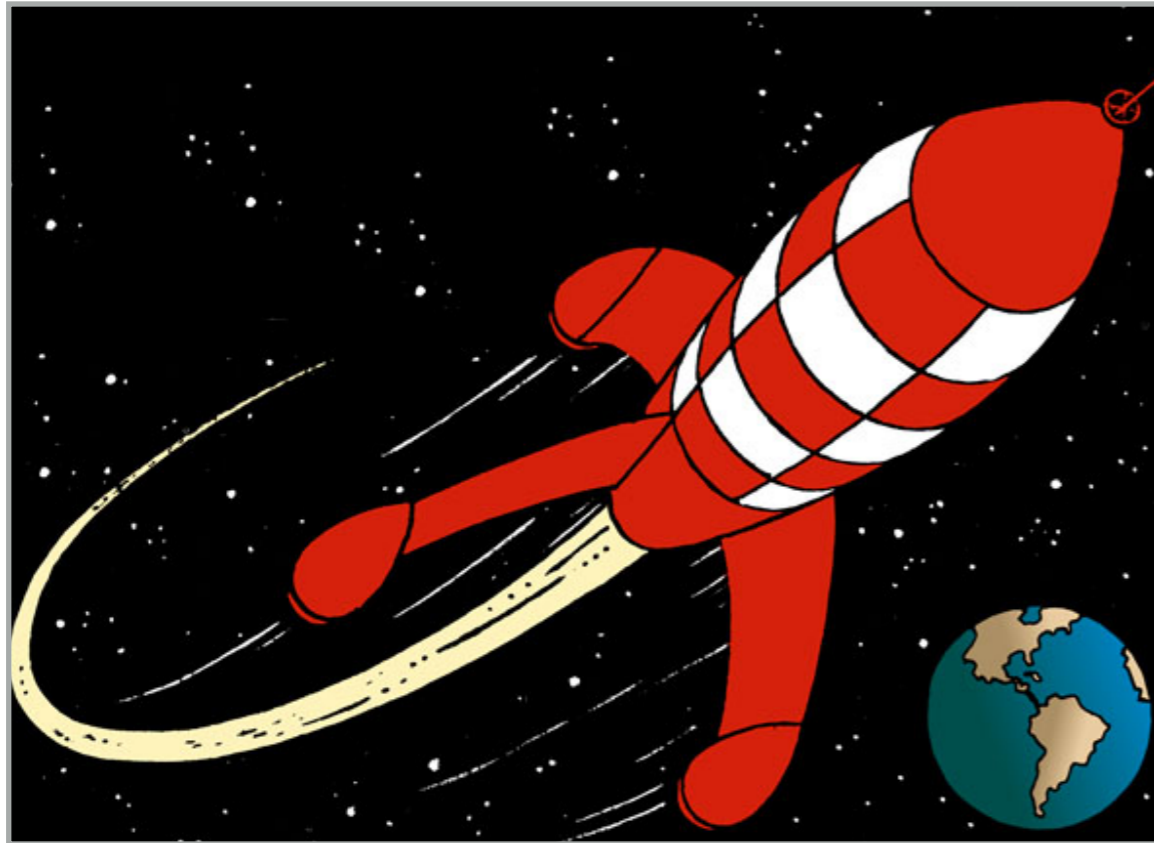


SWIFT: Predictive Fast Reroute



Thomas Holterbach

ETH Zürich / CAIDA

SIGCOMM

24th August 2017

swift.ethz.ch

Joint work with

Stefano Vissicchio

Alberto Dainotti

Laurent Vanbever

UCLondon

CAIDA, UC San Diego

ETH Zürich

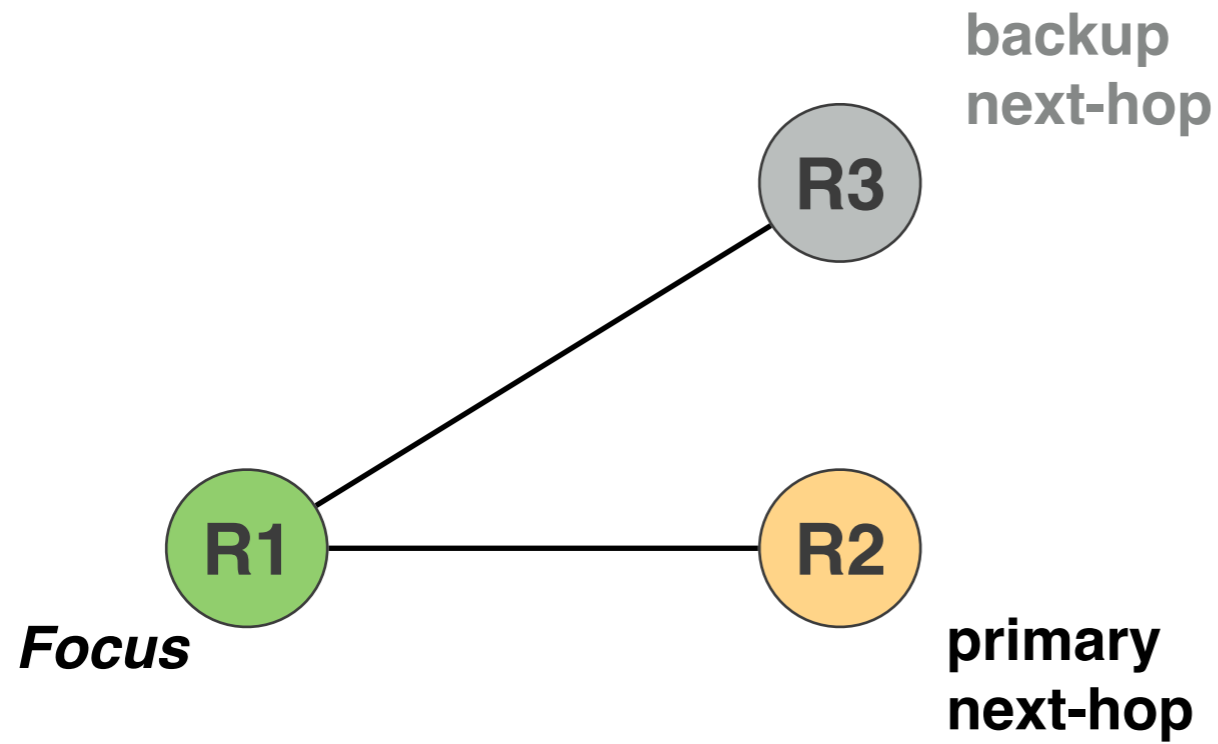
2.5 min

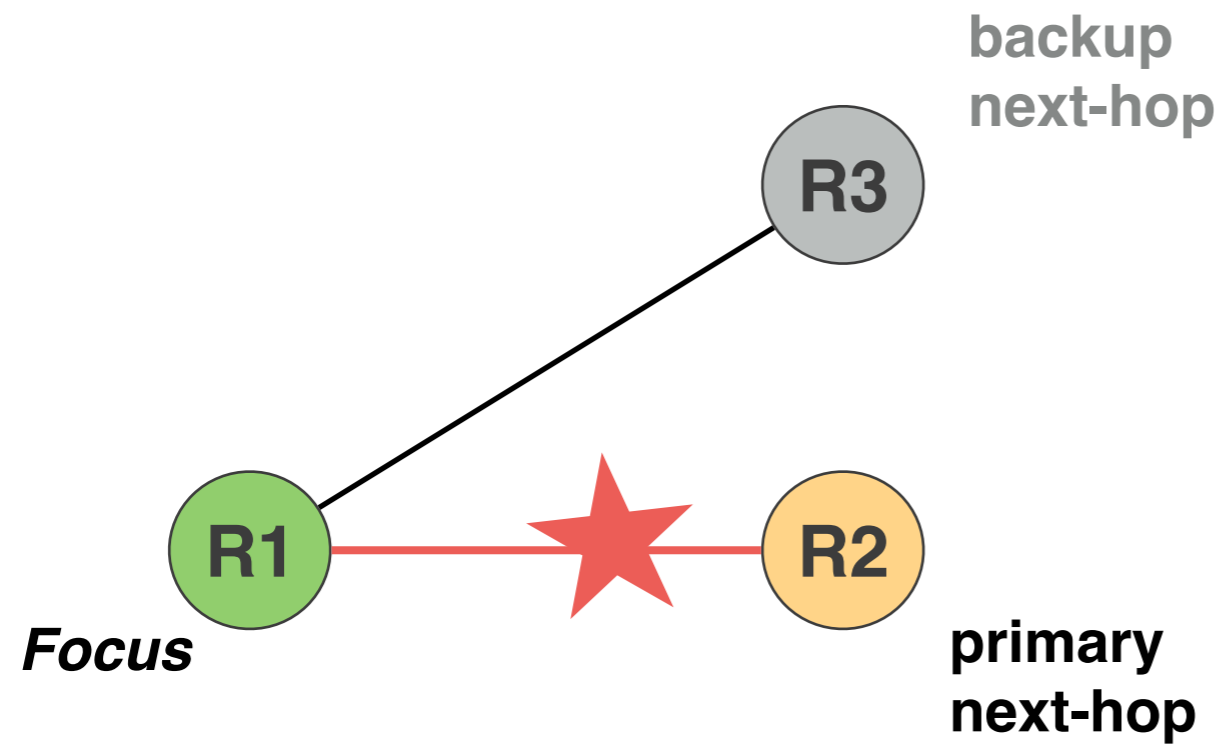
2.5 min



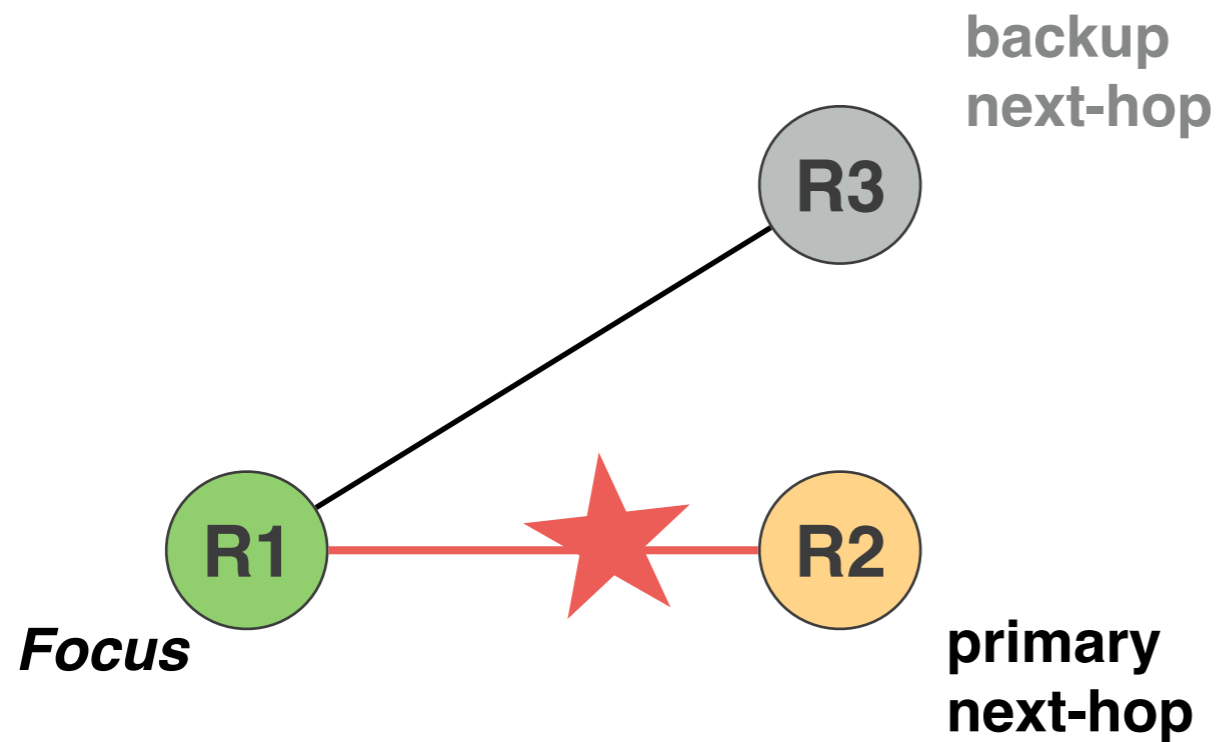
**worst-case convergence time
of a router upon an Internet failure**

Let's consider first a simple example

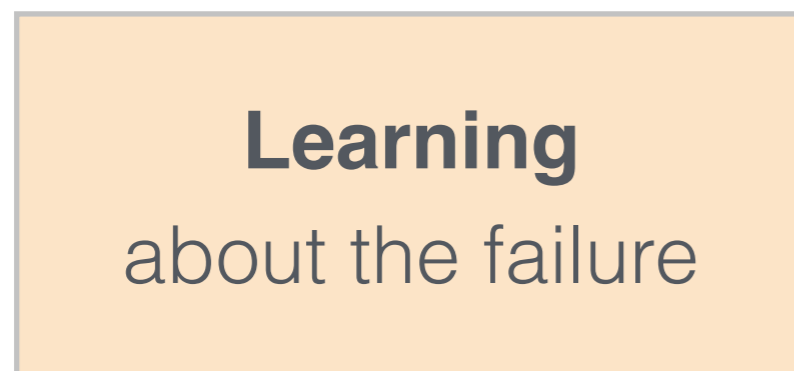




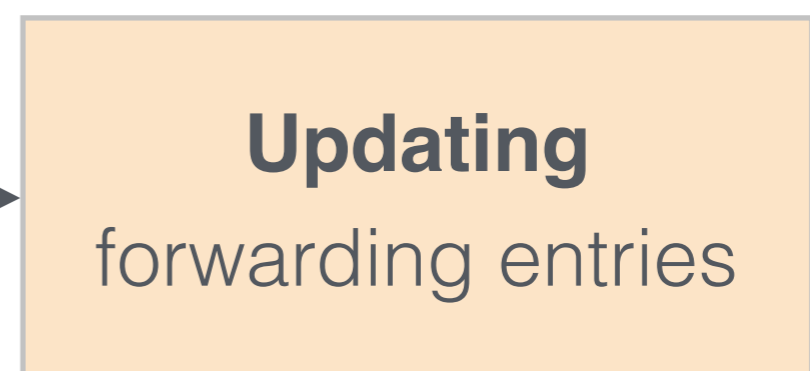
BGP convergence is a two-phase process



Phase #1



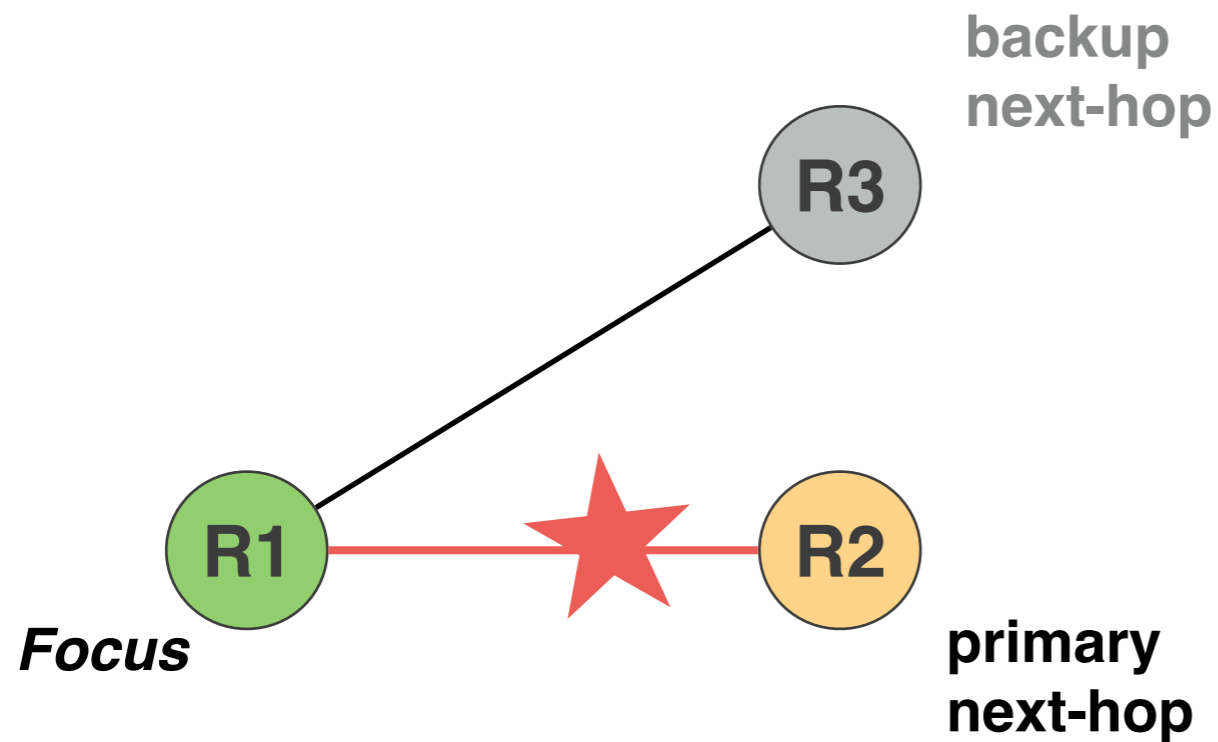
Phase #2



BGP convergence is a two-phase process

R1's forwarding table

pos.	prefix	NH
1	1.0.0.0/24	2
2	1.0.1.0/24	2
3	1.0.2.0/24	2
...
300K	100.0.0.0/24	2
...
600K	200.0.0.0/24	2



Phase #1

Learning
about the failure

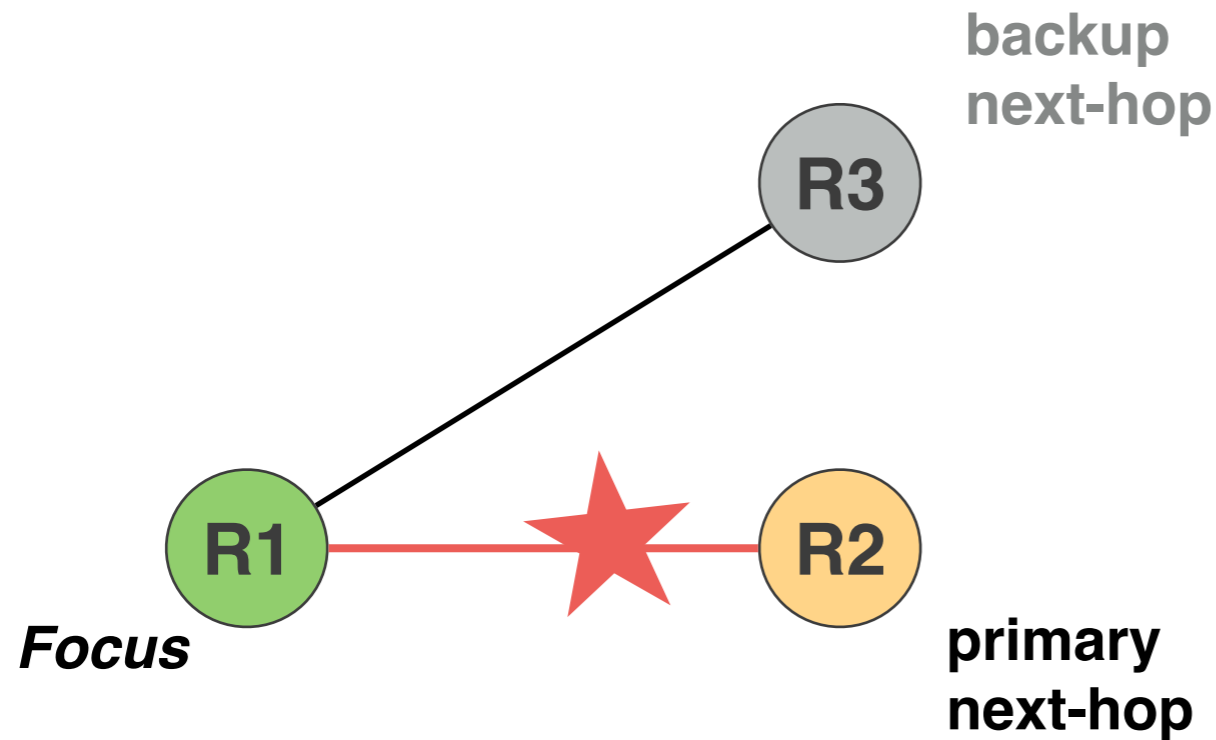
Phase #2

Updating
forwarding entries

BGP convergence is a two-phase process

R1's forwarding table

pos.	prefix	NH
1	1.0.0.0/24	3
2	1.0.1.0/24	2
3	1.0.2.0/24	2
...
300K	100.0.0.0/24	2
...
600K	200.0.0.0/24	2



Phase #1

Learning
about the failure

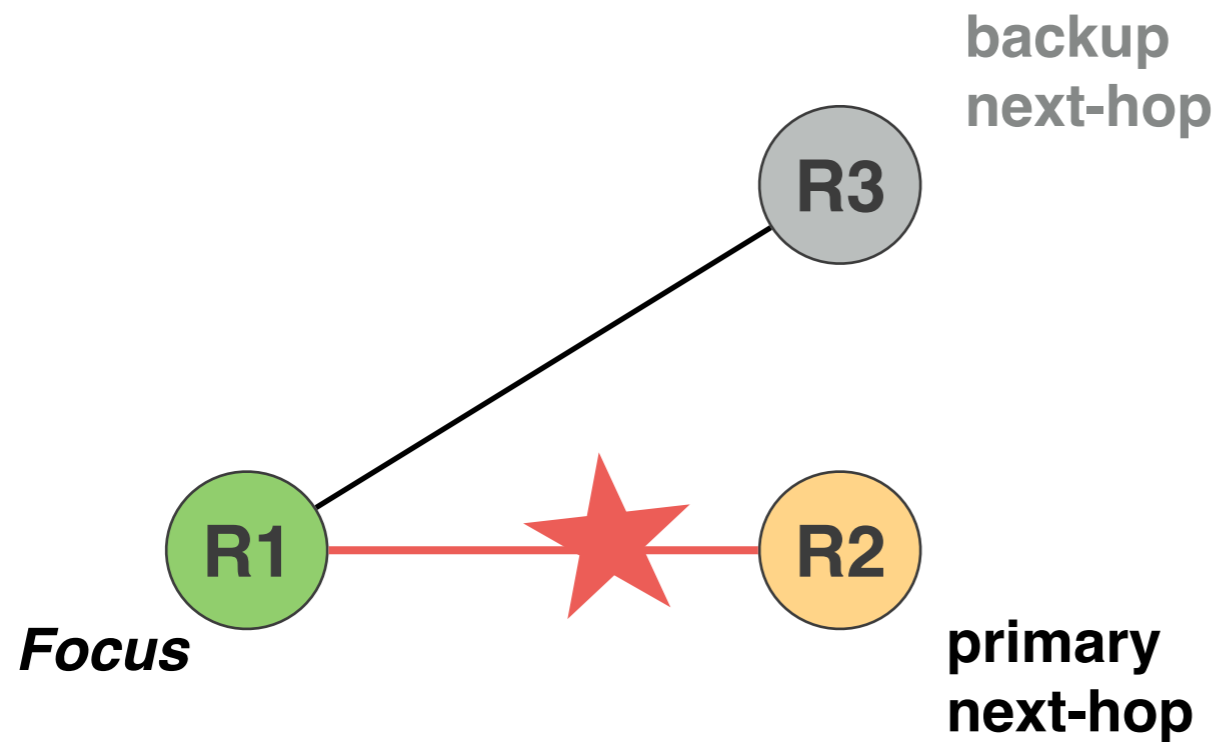
Phase #2

Updating
forwarding entries

BGP convergence is a two-phase process

R1's forwarding table

pos.	prefix	NH
1	1.0.0.0/24	3
2	1.0.1.0/24	3
3	1.0.2.0/24	2
...
300K	100.0.0.0/24	2
...
600K	200.0.0.0/24	2



Phase #1

Learning
about the failure

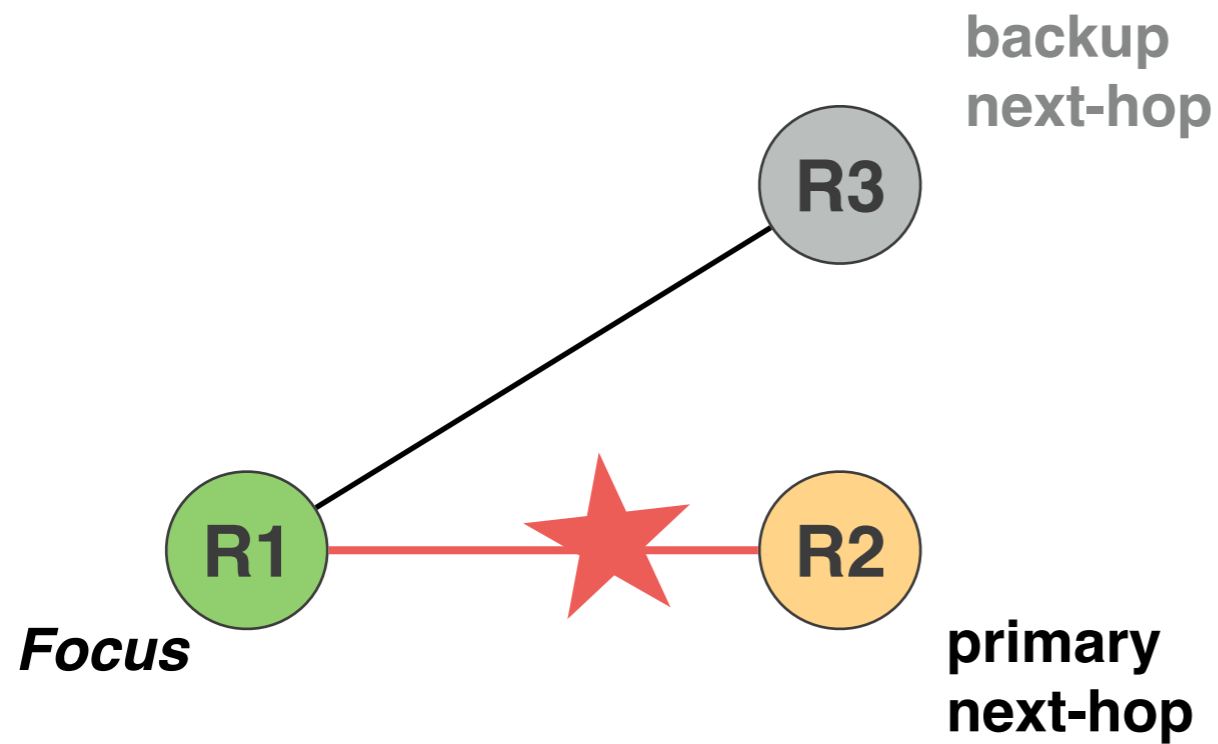
Phase #2

Updating
forwarding entries

BGP convergence is a two-phase process

R1's forwarding table

pos.	prefix	NH
1	1.0.0.0/24	3
2	1.0.1.0/24	3
3	1.0.2.0/24	3
...
300K	100.0.0.0/24	2
...
600K	200.0.0.0/24	2



Phase #1

Learning
about the failure

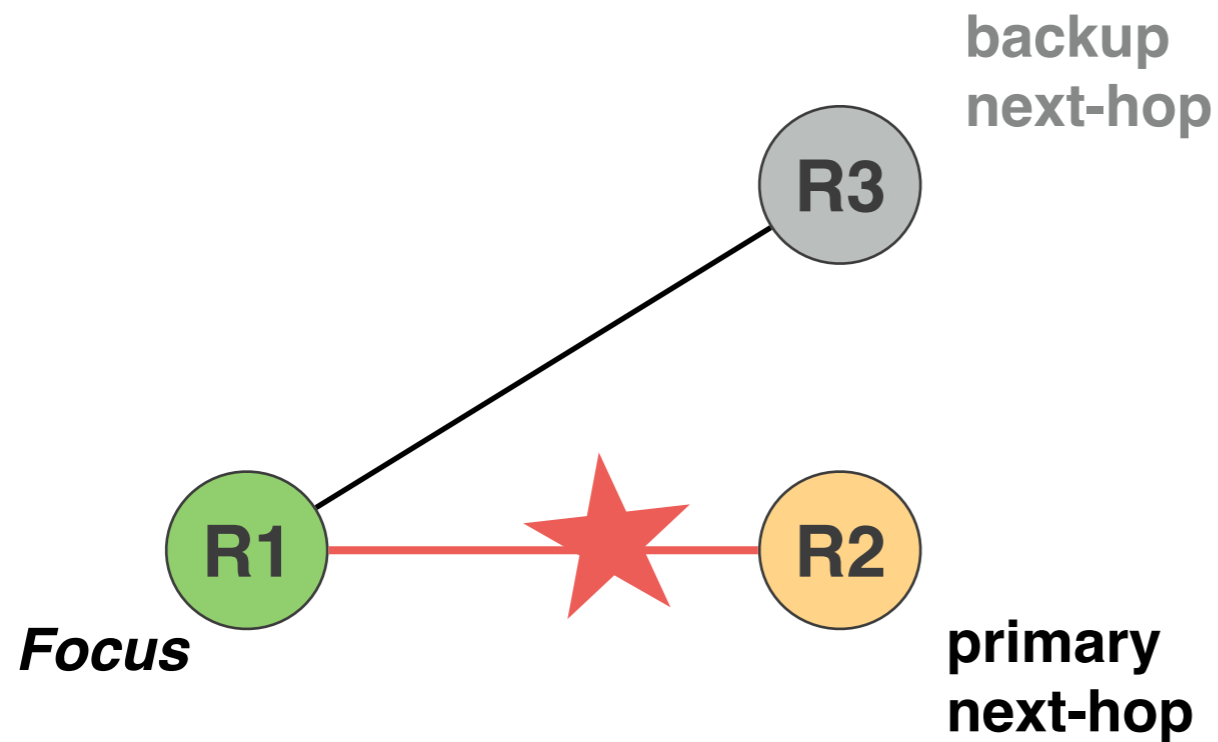
Phase #2

Updating
forwarding entries

BGP convergence is a two-phase process

R1's forwarding table

pos.	prefix	NH
1	1.0.0.0/24	3
2	1.0.1.0/24	3
3	1.0.2.0/24	3
...
300K	100.0.0.0/24	3
...
600K	200.0.0.0/24	2



Phase #1

Learning
about the failure

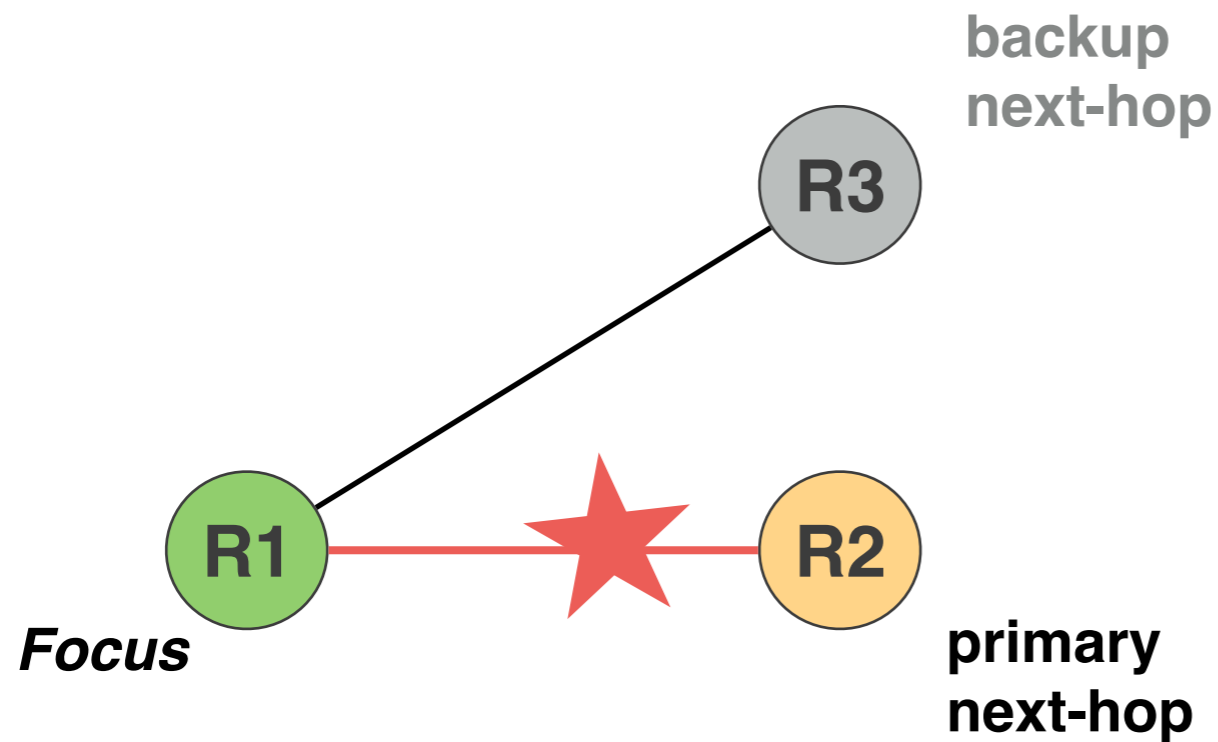
Phase #2

Updating
forwarding entries

BGP convergence is a two-phase process

R1's forwarding table

pos.	prefix	NH
1	1.0.0.0/24	3
2	1.0.1.0/24	3
3	1.0.2.0/24	3
...
300K	100.0.0.0/24	3
...
600K	200.0.0.0/24	3



Phase #1

Learning
about the failure

Phase #2

Updating
forwarding entries



BGP Fast Reroute

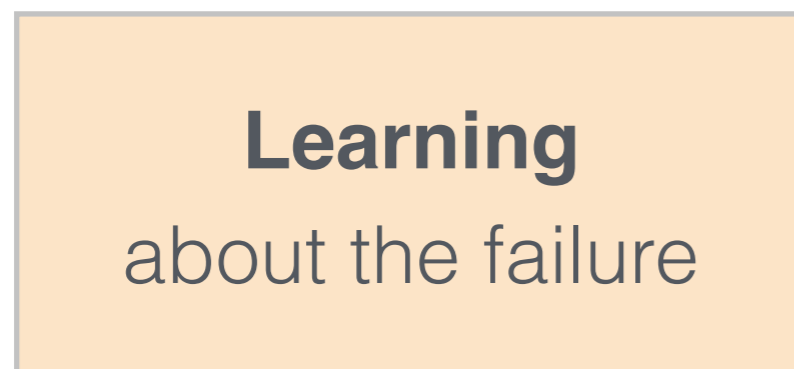
About 3,710 results (0.08 sec)

In practice, multiple techniques are often used to speed up the convergence time

BGP timers
BFD, ...



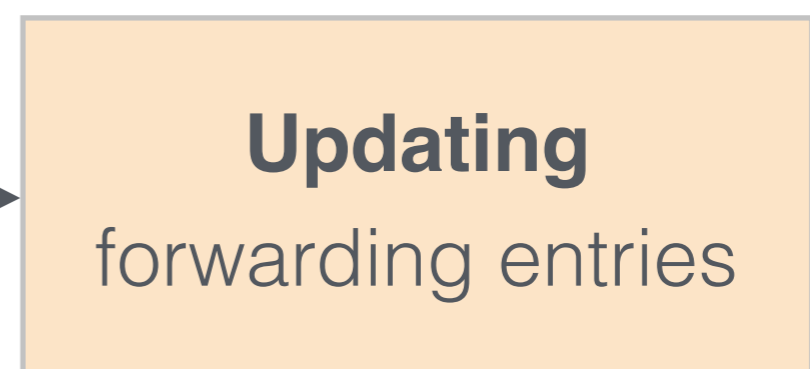
Phase #1



PIC
MPLS, ...



Phase #2



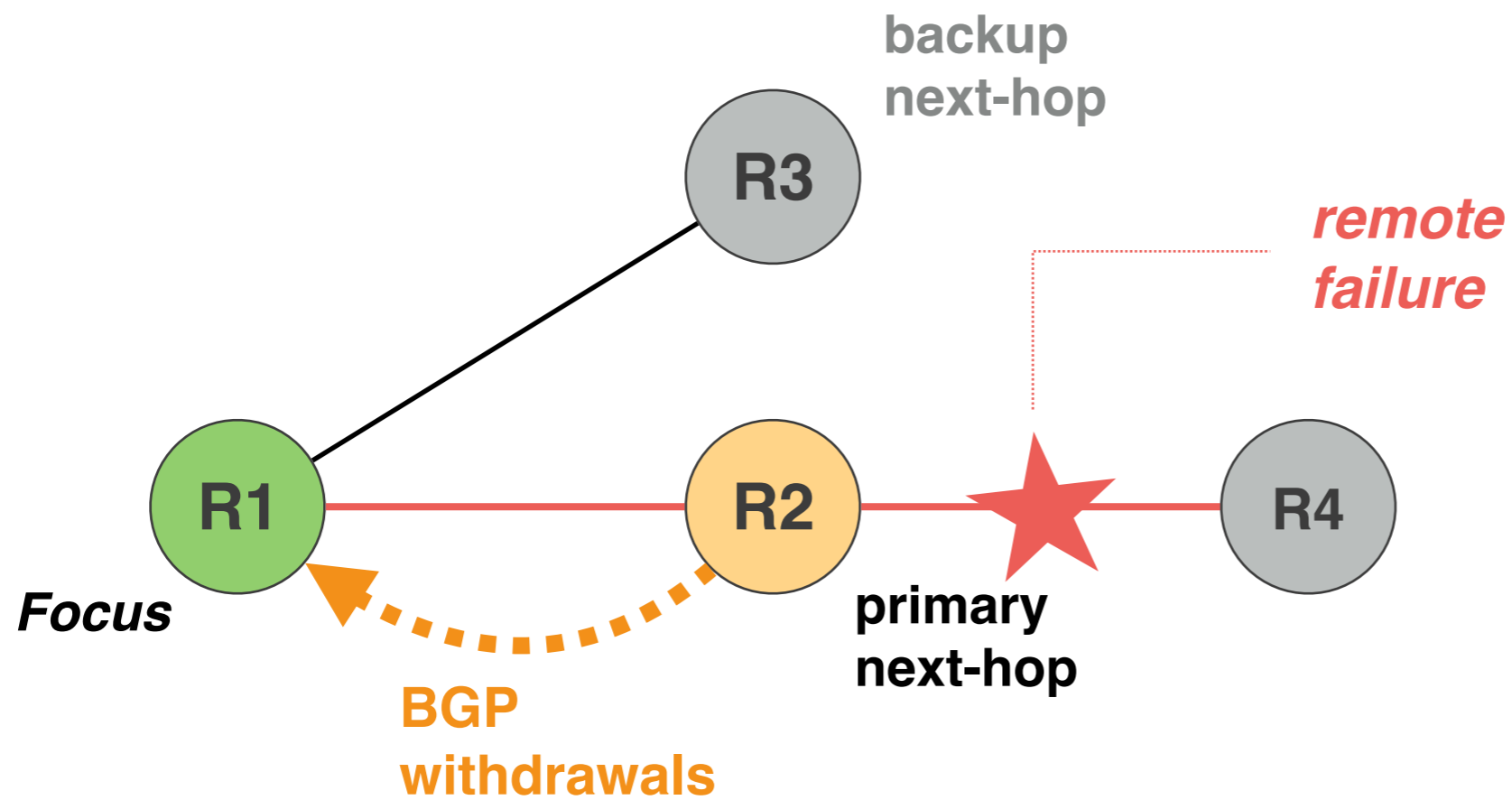


BGP Fast Reroute

About 3,710 results (0.08 sec)

**none of them work
upon remote outages**

Current fast reroute techniques do not work upon **remote outages**



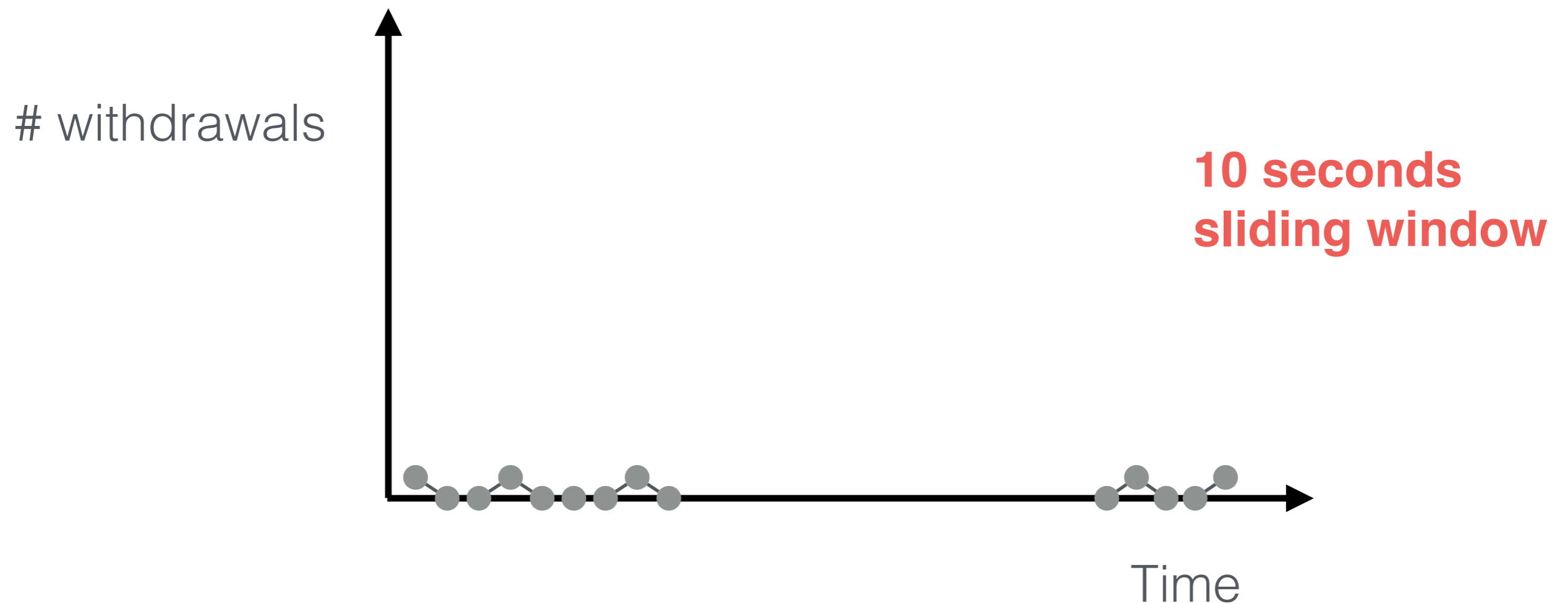
Remote outages are numerous

Remote outages **are numerous**

We analyzed BGP sessions from RouteViews and RIPE RIS collectors, and looked for **bursts of BGP withdrawals**

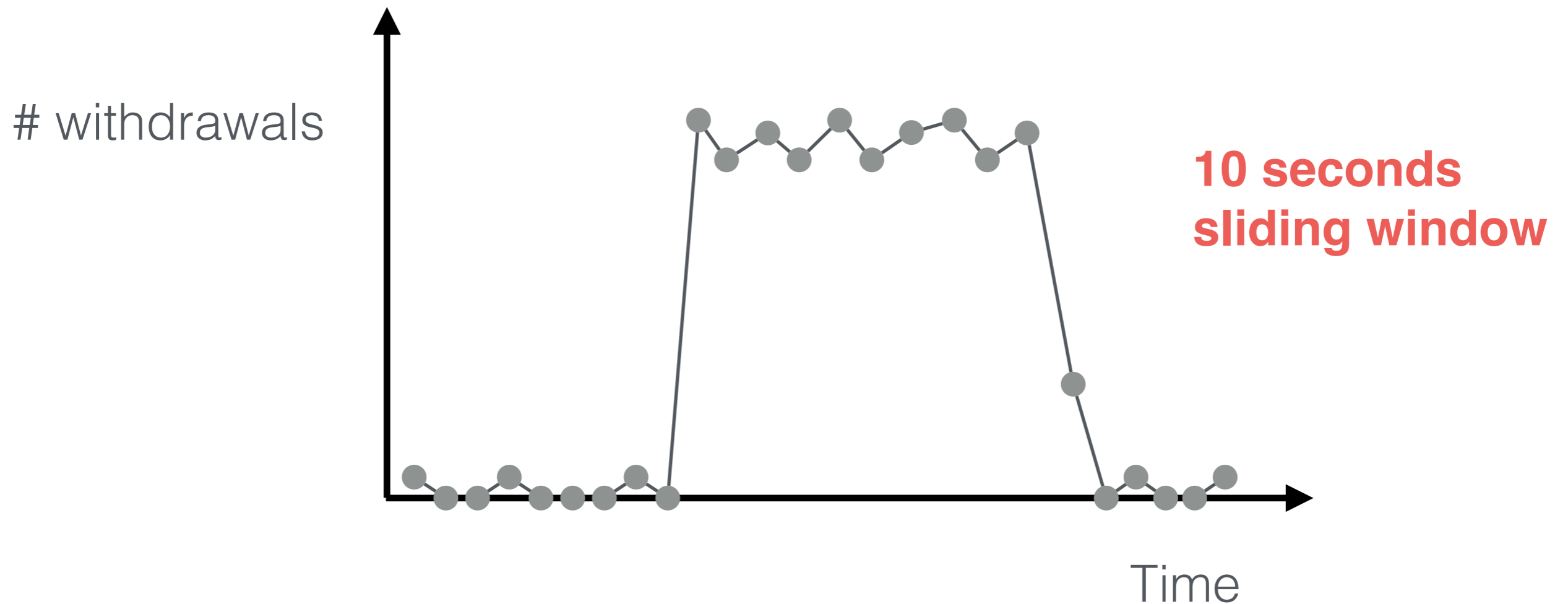
Remote outages **are numerous**

We analyzed BGP sessions from RouteViews and RIPE RIS collectors, and looked for **bursts of BGP withdrawals**



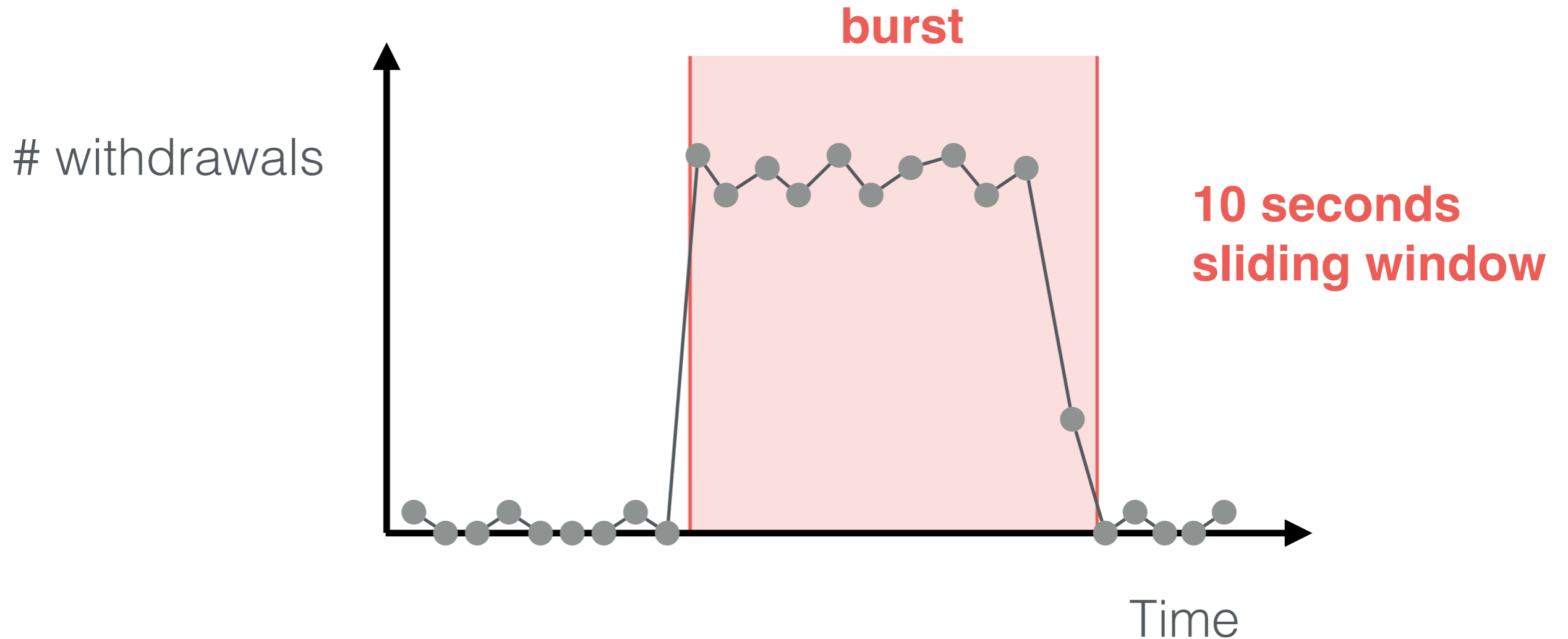
Remote outages **are numerous**

We analyzed BGP sessions from RouteViews and RIPE RIS collectors, and looked for **bursts of BGP withdrawals**



Remote outages **are numerous**

We analyzed BGP sessions from RouteViews and RIPE RIS collectors, and looked for **bursts of BGP withdrawals**

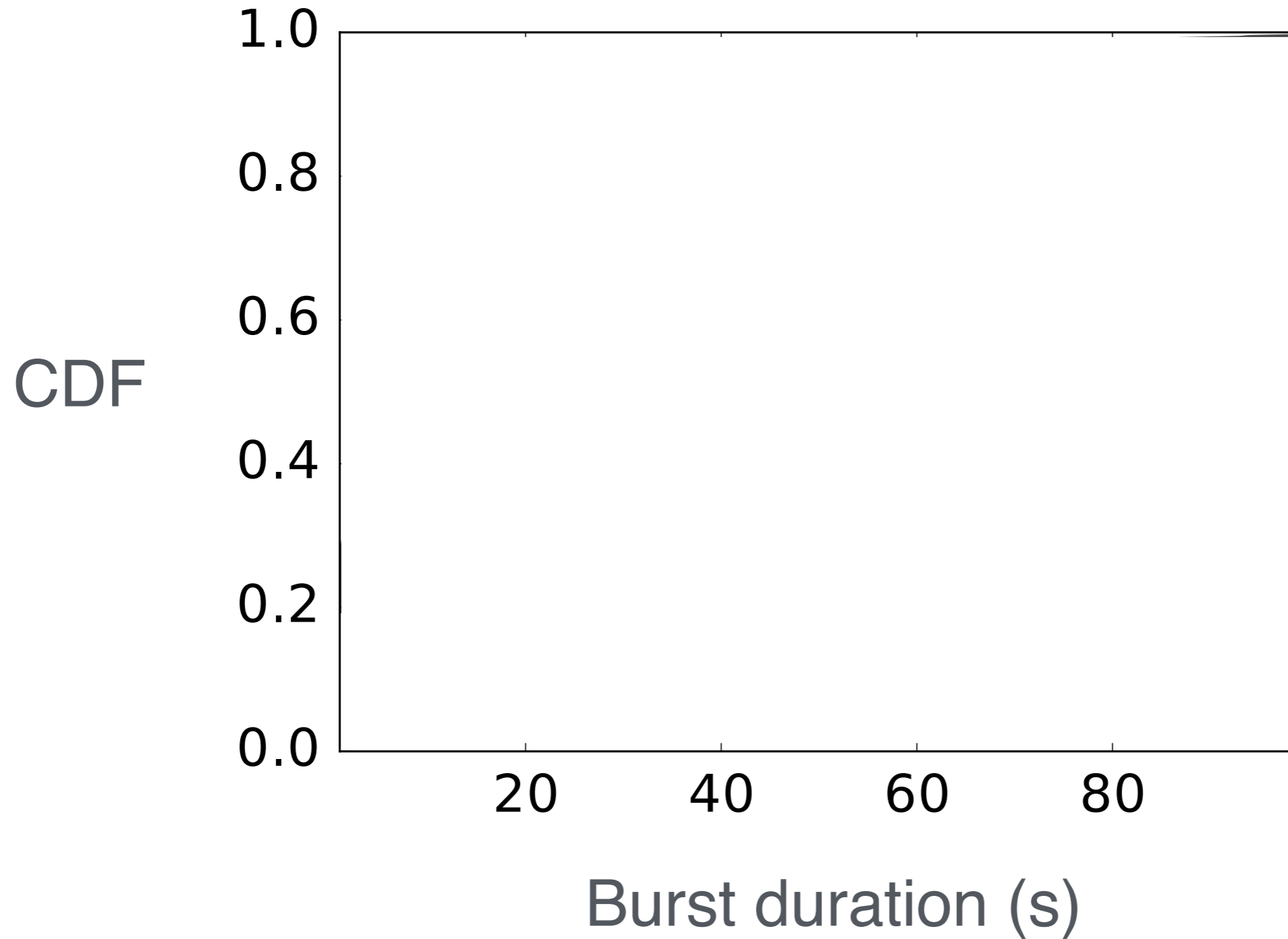


Remote outages **are numerous**

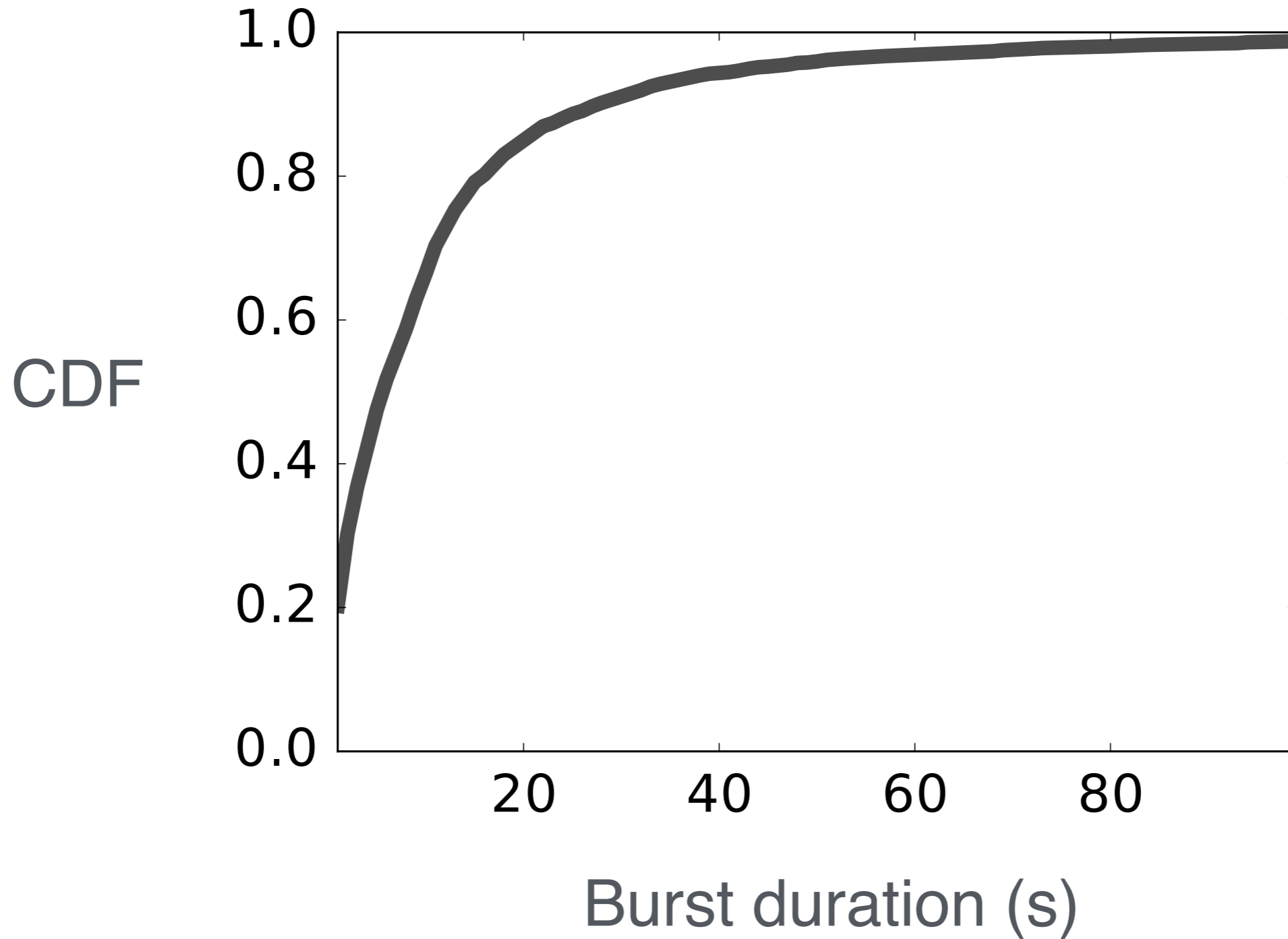
We analyzed BGP sessions from RouteViews and RIPE RIS collectors, and looked for **bursts of BGP withdrawals**

We found **~3,000 bursts** of withdrawals in November 2016 and on 213 BGP sessions

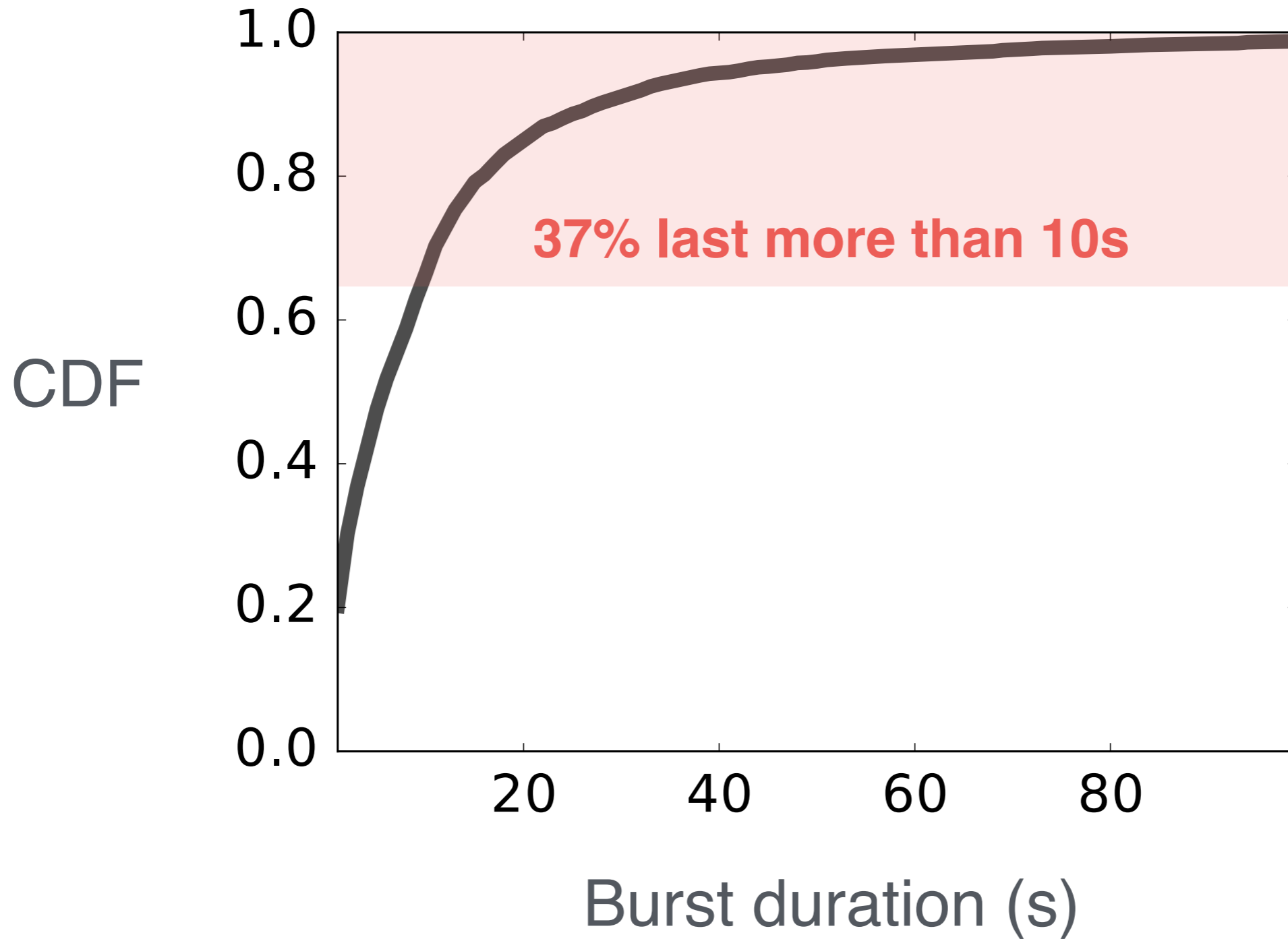
Remote outages propagate slowly



Remote outages propagate slowly

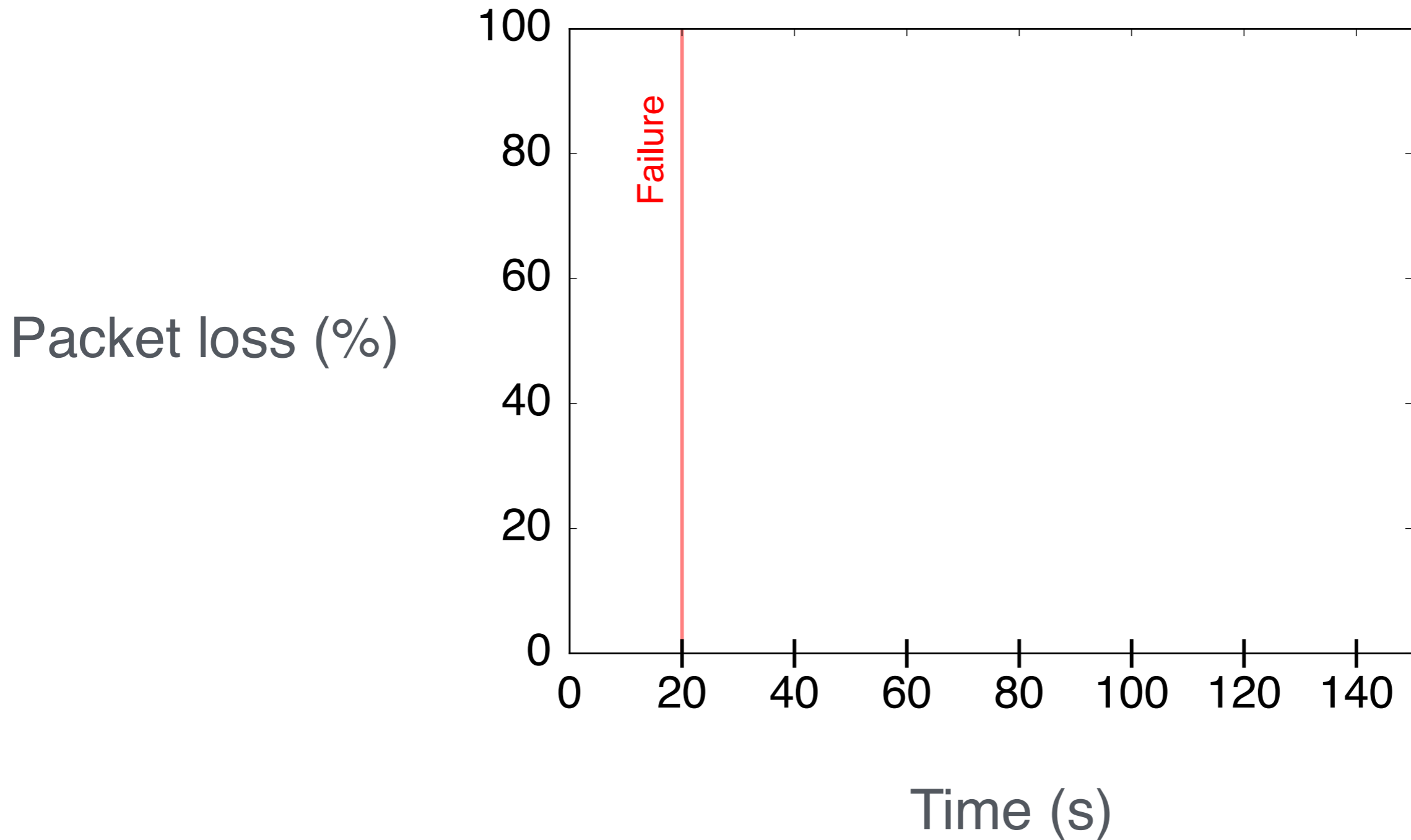


Remote outages propagate slowly

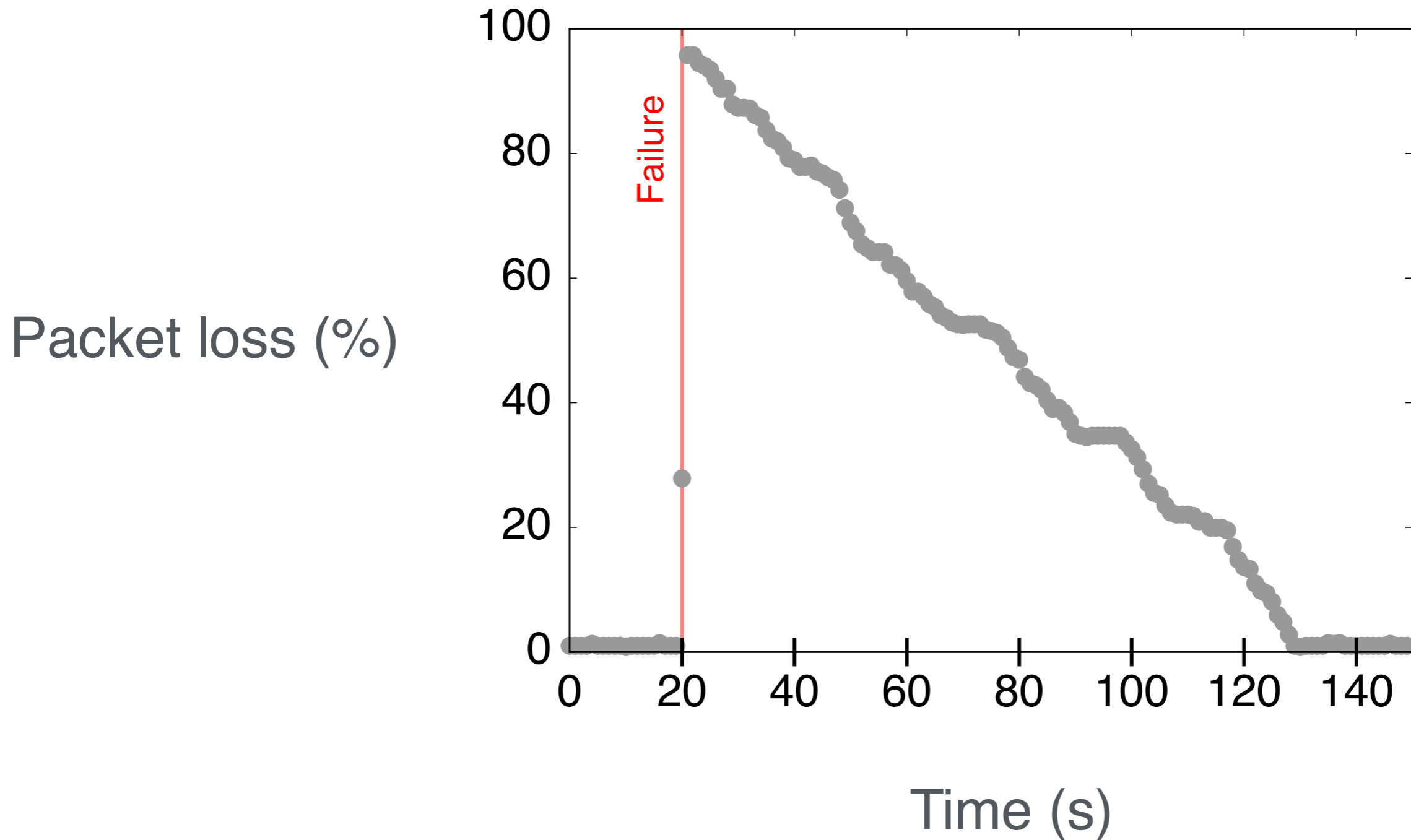


Remote outages **result in downtime**

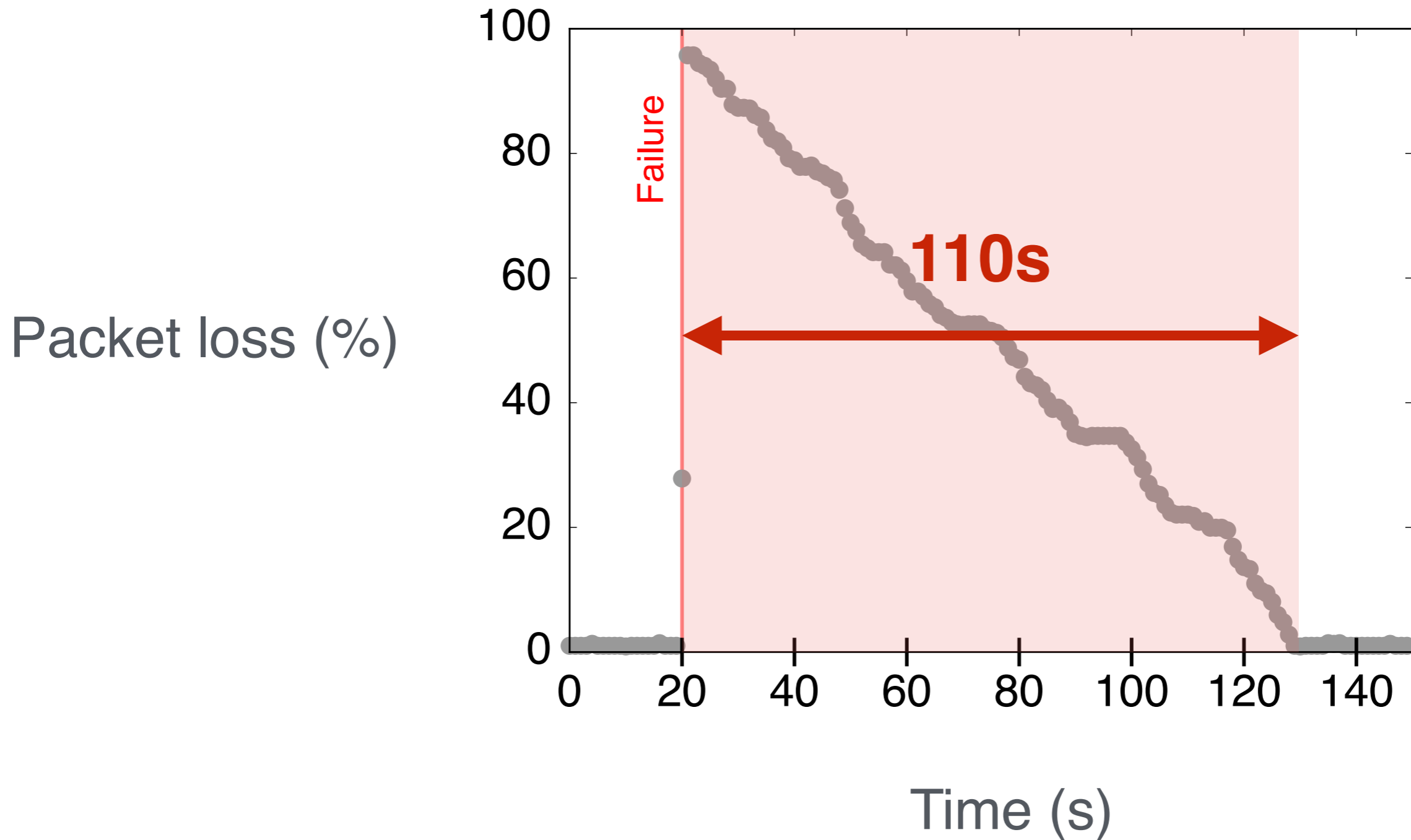
Remote outages result in downtime



Remote outages result in downtime



Remote outages result in downtime



SWIFT: Predictive Fast Reroute Framework

SWIFT: Predictive Fast Reroute Framework

works upon **remote** (and local) outages

SWIFT: Predictive Fast Reroute

1. ***SWIFT*** reduces the learning time of a withdrawal from **13s to 2s** using a control-plane prediction algorithm

SWIFT: Predictive Fast Reroute

1. **SWIFT** reduces the learning time of a withdrawal from **13s to 2s** using a control-plane prediction algorithm
2. **SWIFT** fast reroutes the affected traffic towards unaffected paths in **40ms** using a hierarchical forwarding table

SWIFT: Predictive Fast Reroute

1. **SWIFT** reduces the learning time of a withdrawal from **13s to 2s** using a control-plane prediction algorithm
2. **SWIFT** fast reroutes the affected traffic towards unaffected paths in **40ms** using a hierarchical forwarding table
3. **SWIFT** is **deployable** on existing routers

SWIFT: Predictive Fast Reroute

1. **SWIFT** reduces the learning time of a withdrawal from **13s to 2s** using a control-plane prediction algorithm
2. **SWIFT** fast reroutes the affected traffic towards unaffected paths in **40ms** using a hierarchical forwarding table
3. **SWIFT** is **deployable** on existing routers

SWIFT predicts the extent of an outage
using **fast, yet precise** Root Cause Analysis

Existing RCA techniques *

accuracy

vantage points

detection speed

- * Lifeguard, SIGCOMM'14
- Poiroot, SIGCOMM'13
- Feldmann et al., SIGCOMM'04

Existing RCA techniques *

accuracy

link/AS

vantage points

multiple

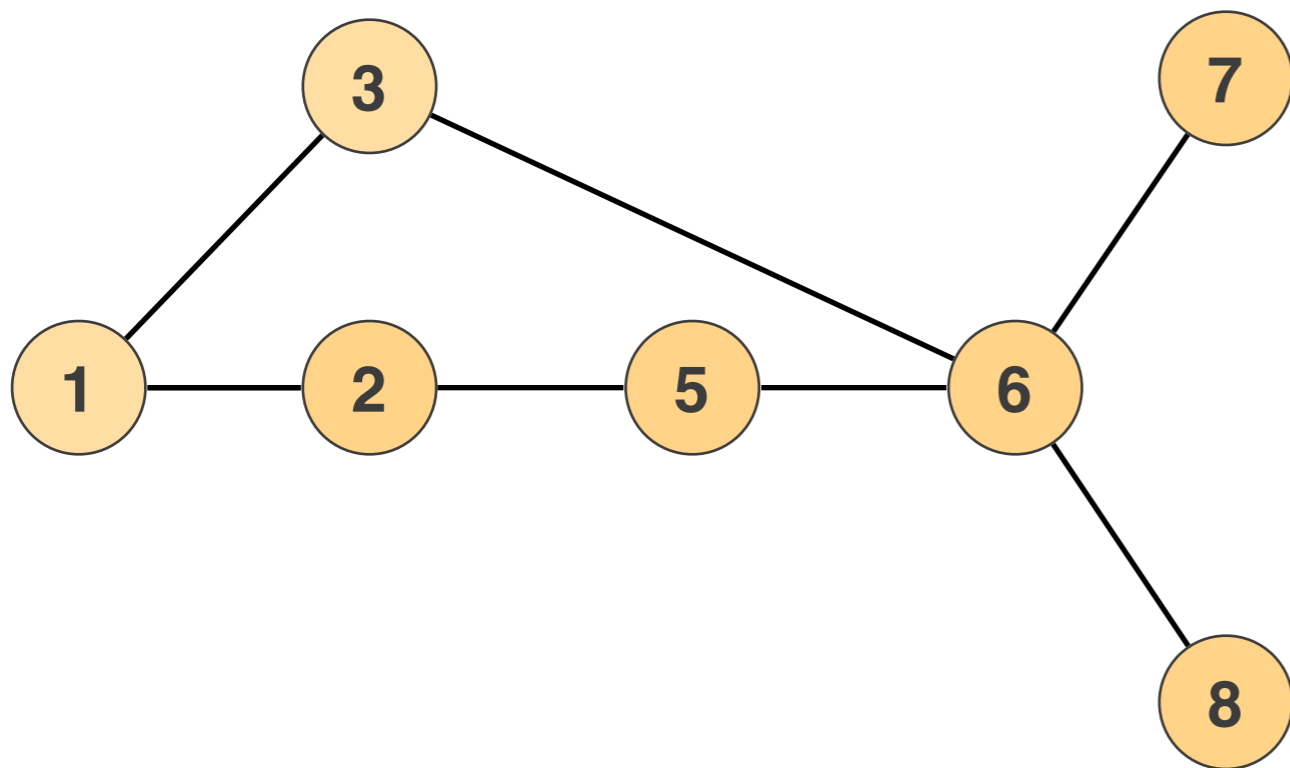
detection speed

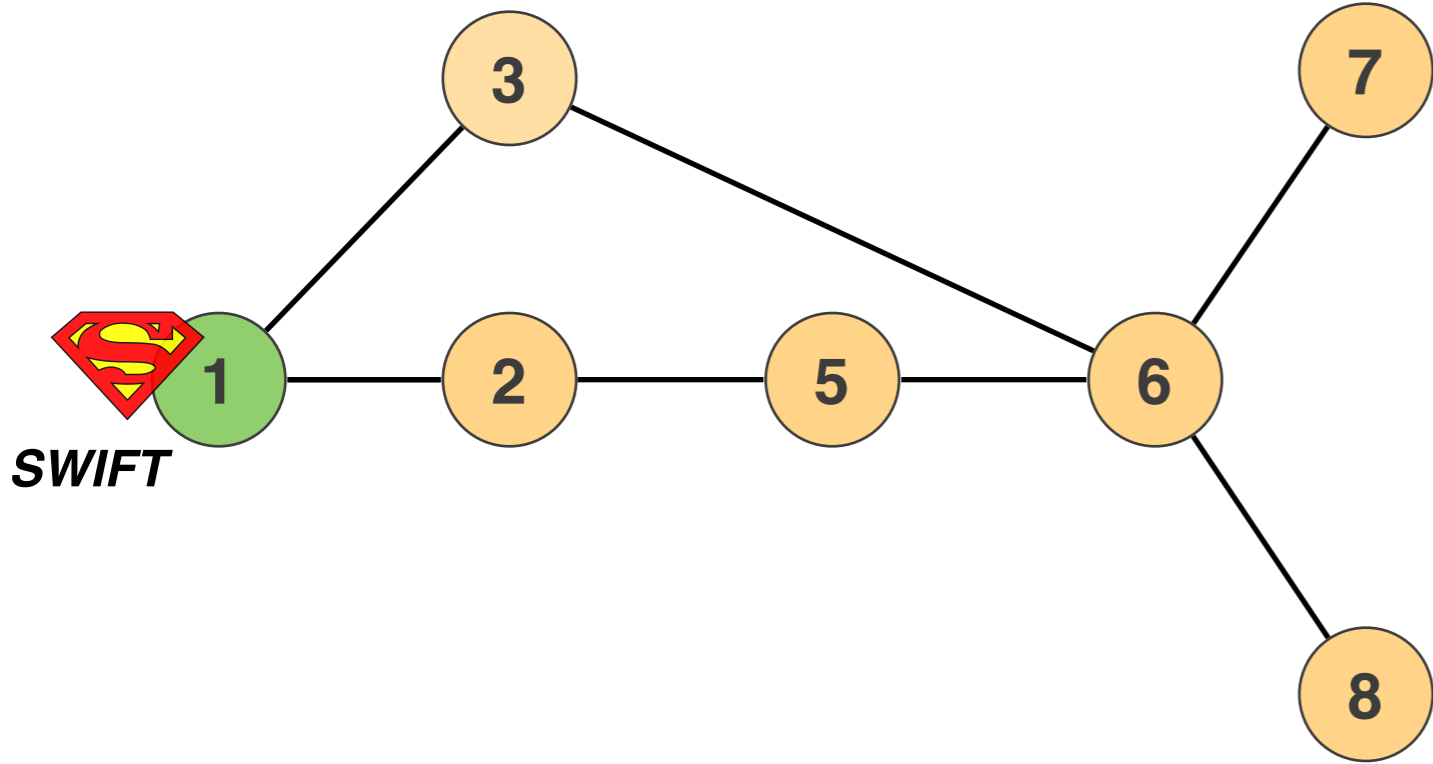
O(minute)

- * Lifeguard, SIGCOMM'14
- Poiroot, SIGCOMM'13
- Feldmann et al., SIGCOMM'04

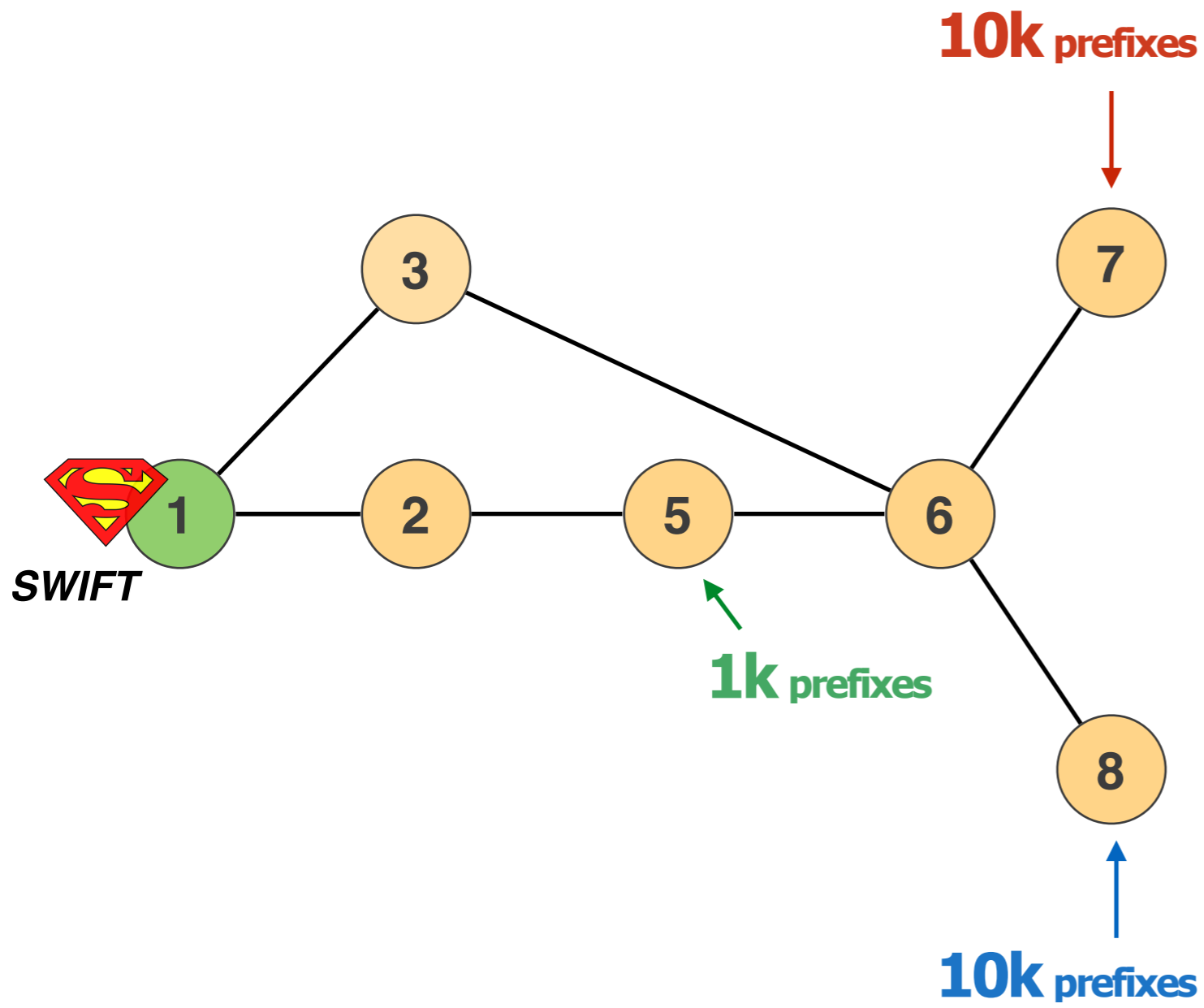
	Existing RCA techniques *	SWIFT
accuracy	link/AS	region
# vantage points	multiple	1
detection speed	O(minute)	O(second)

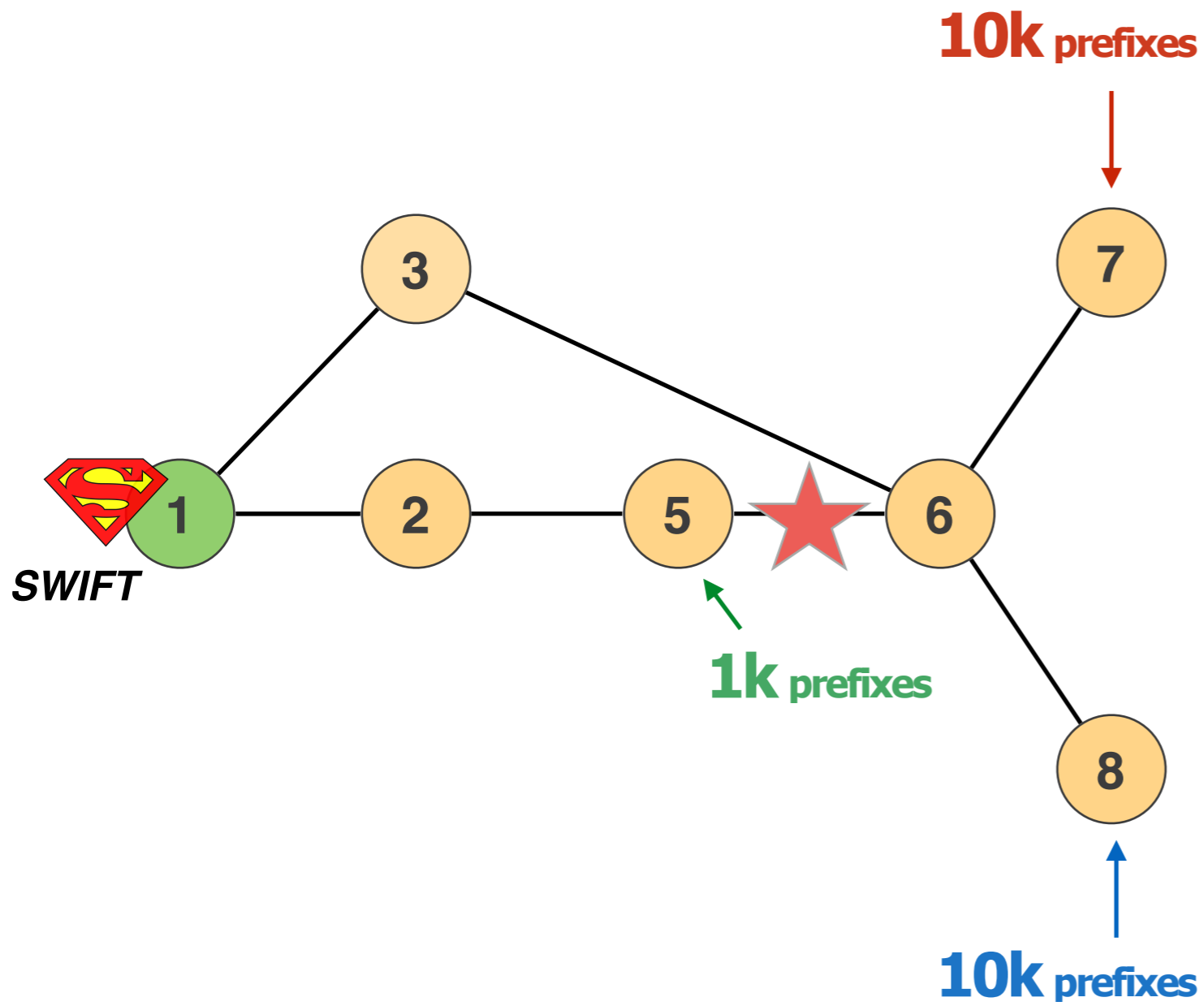
- * Lifeguard, SIGCOMM'14
- Poiroot, SIGCOMM'13
- Feldmann et al., SIGCOMM'04

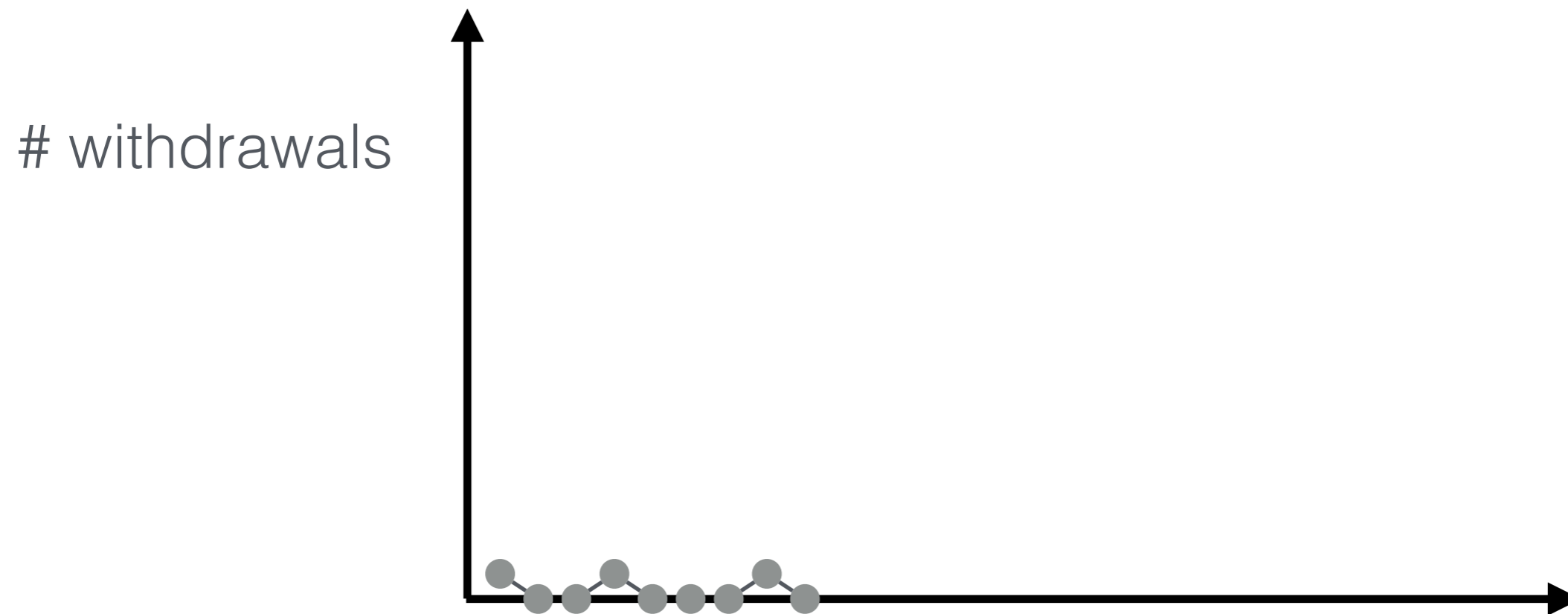




SWIFT







#1 monitor the stream
of BGP messages

Time

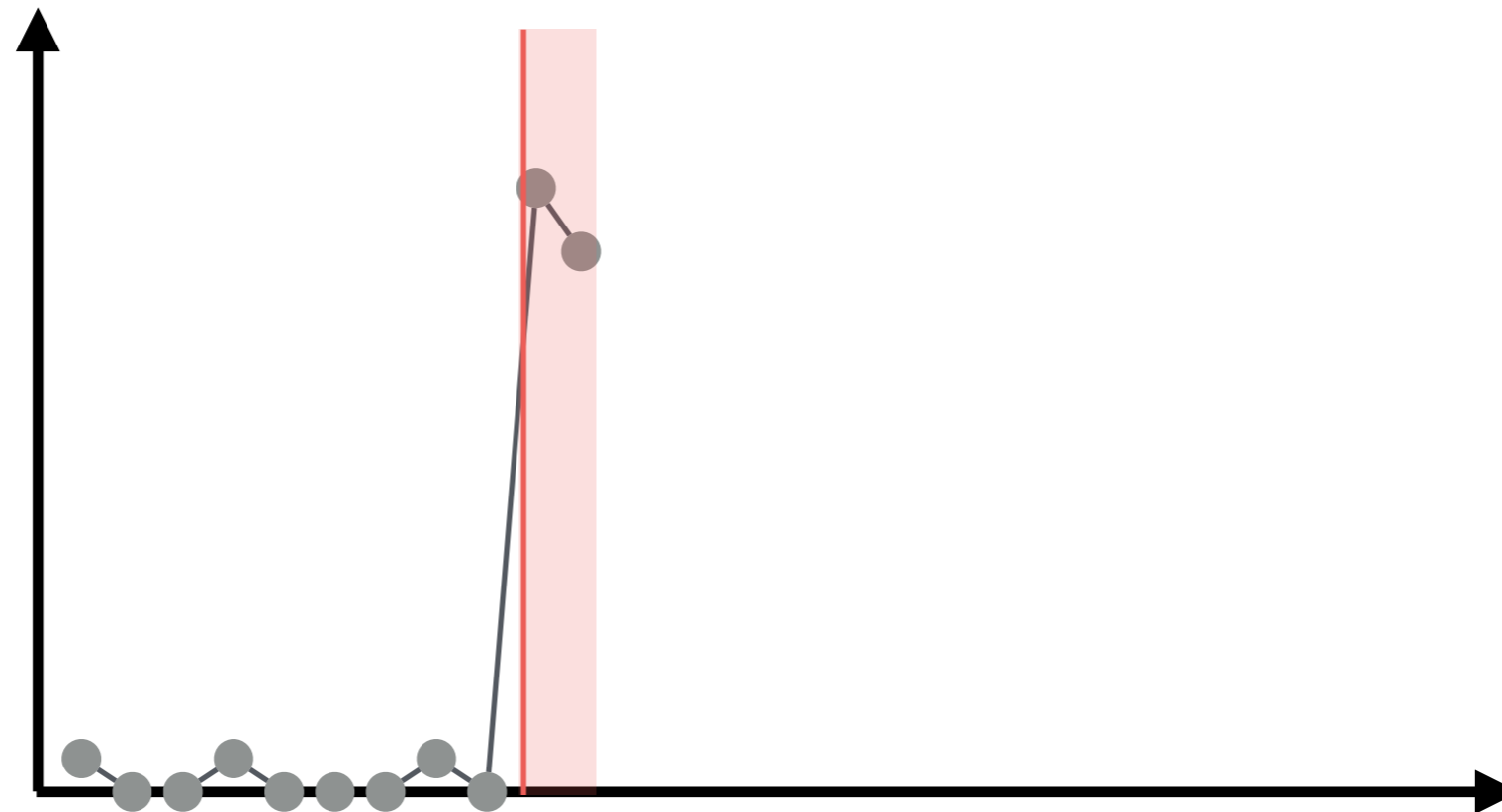
#2 detect and locate the outage

Links failure probability

(A,B)	0.3
⋮	⋮
(A,D)	0.9

→ Link **(A,D)** is dead

withdrawals



#1 monitor the stream of BGP messages

Time

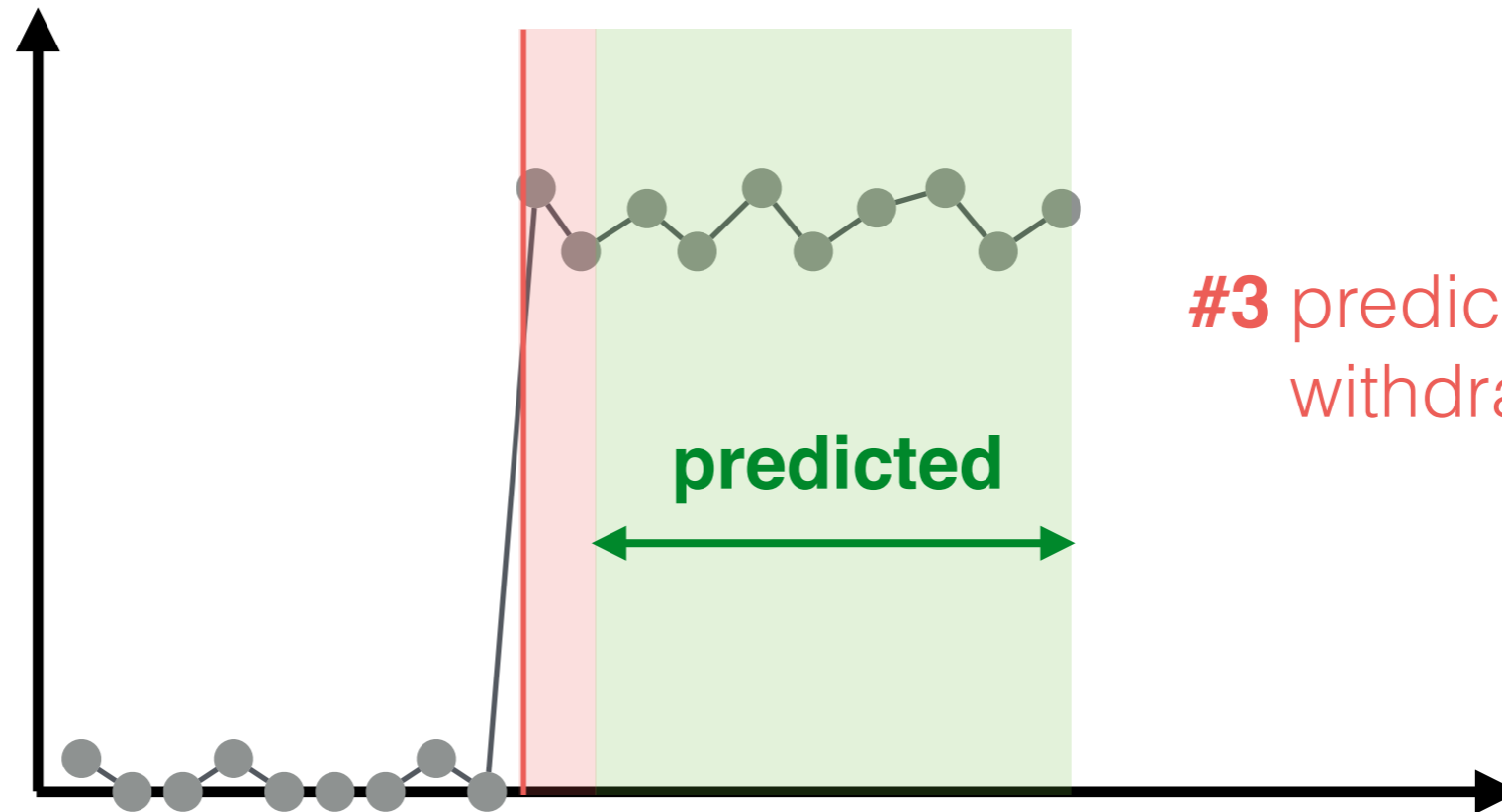
#2 detect and locate the outage

Links failure probability

(A,B)	0.3
⋮	⋮
(A,D)	0.9

→ Link **(A,D)** is dead

withdrawals

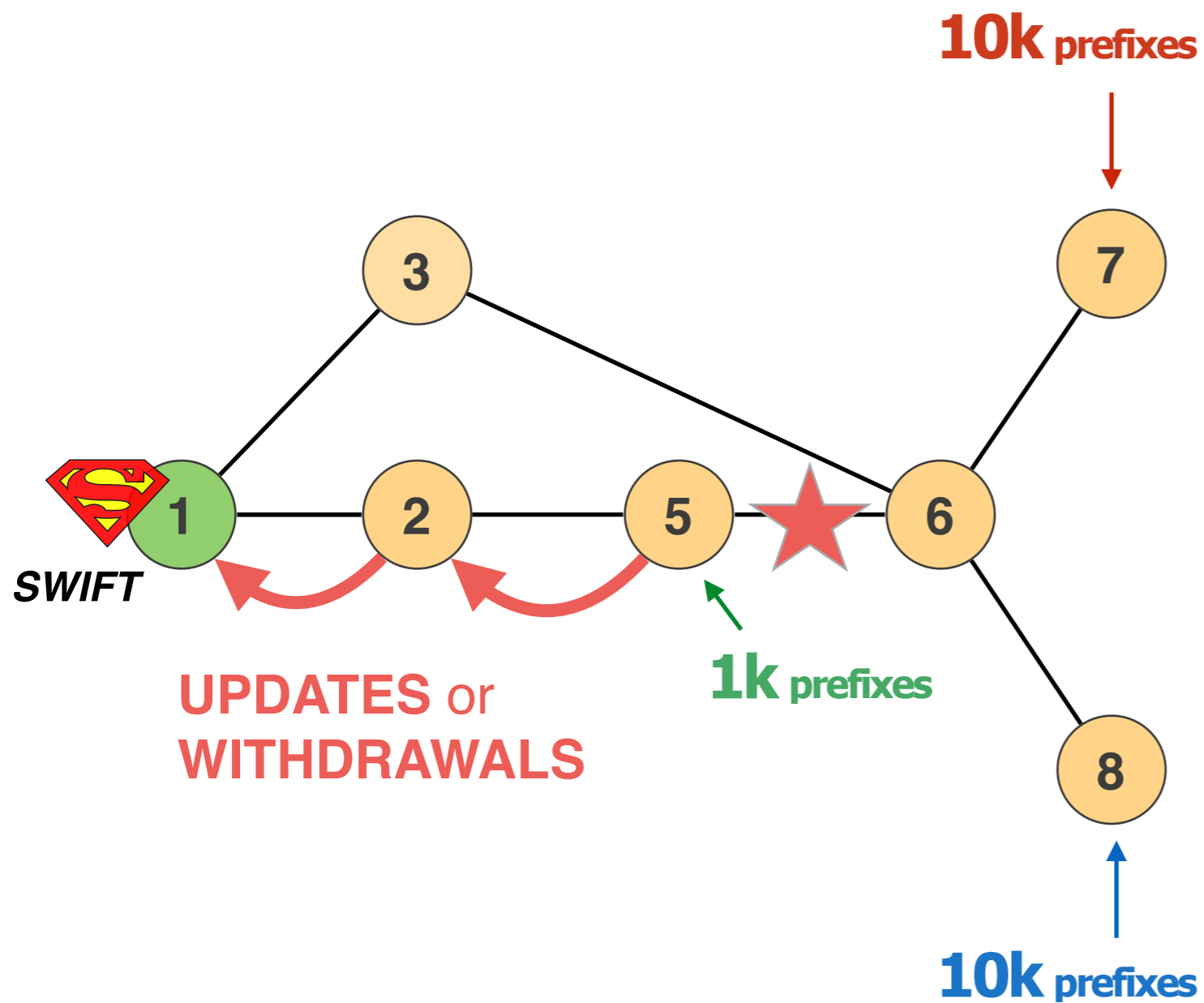


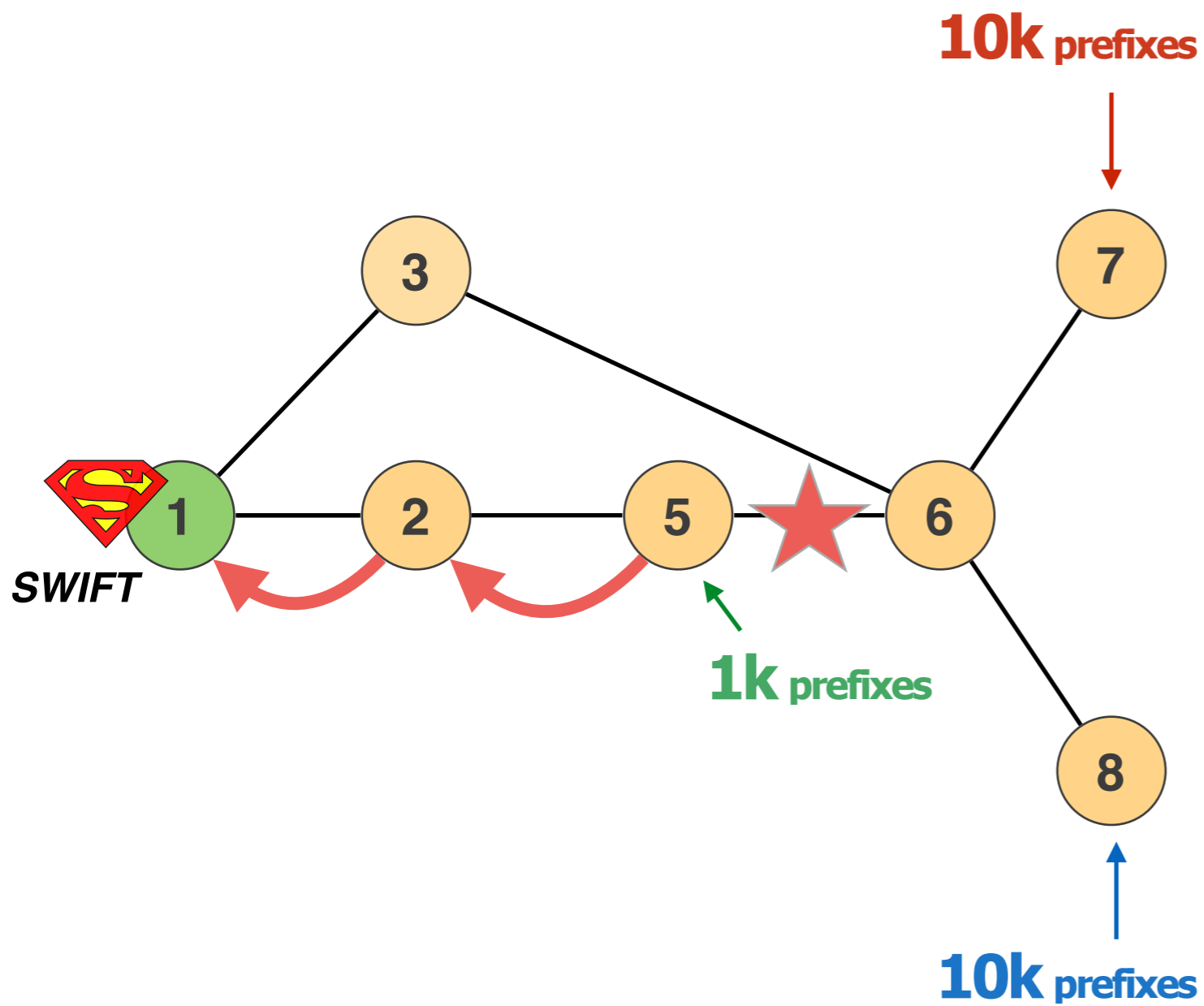
#3 predict future withdrawals

#1 monitor the stream of BGP messages

Time

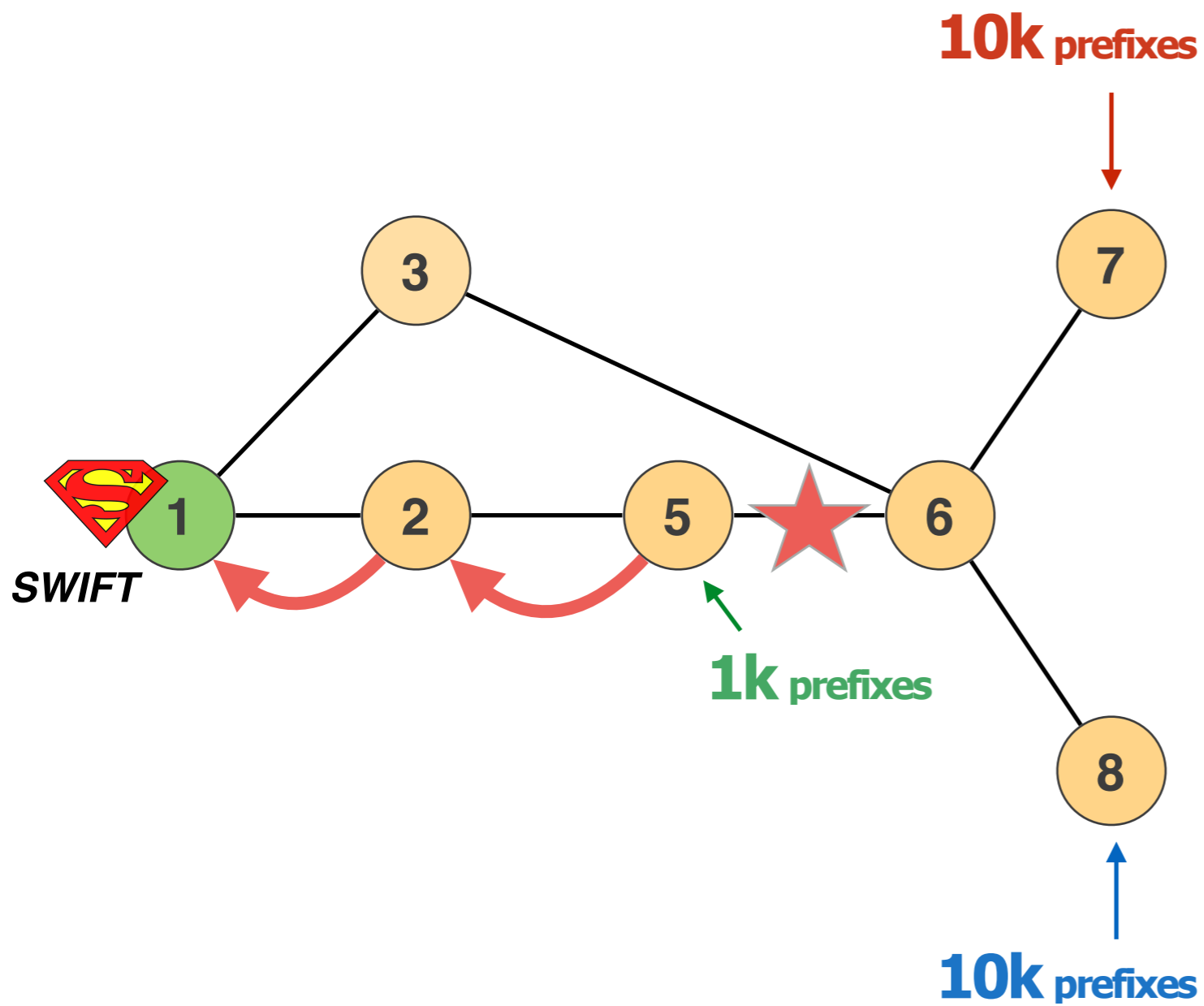
SWIFT link failure inference algorithm leverages the **affected** and **unaffected** prefixes





AS links

candidates set
 (2,5) (5,6) (6,7) (6,8)



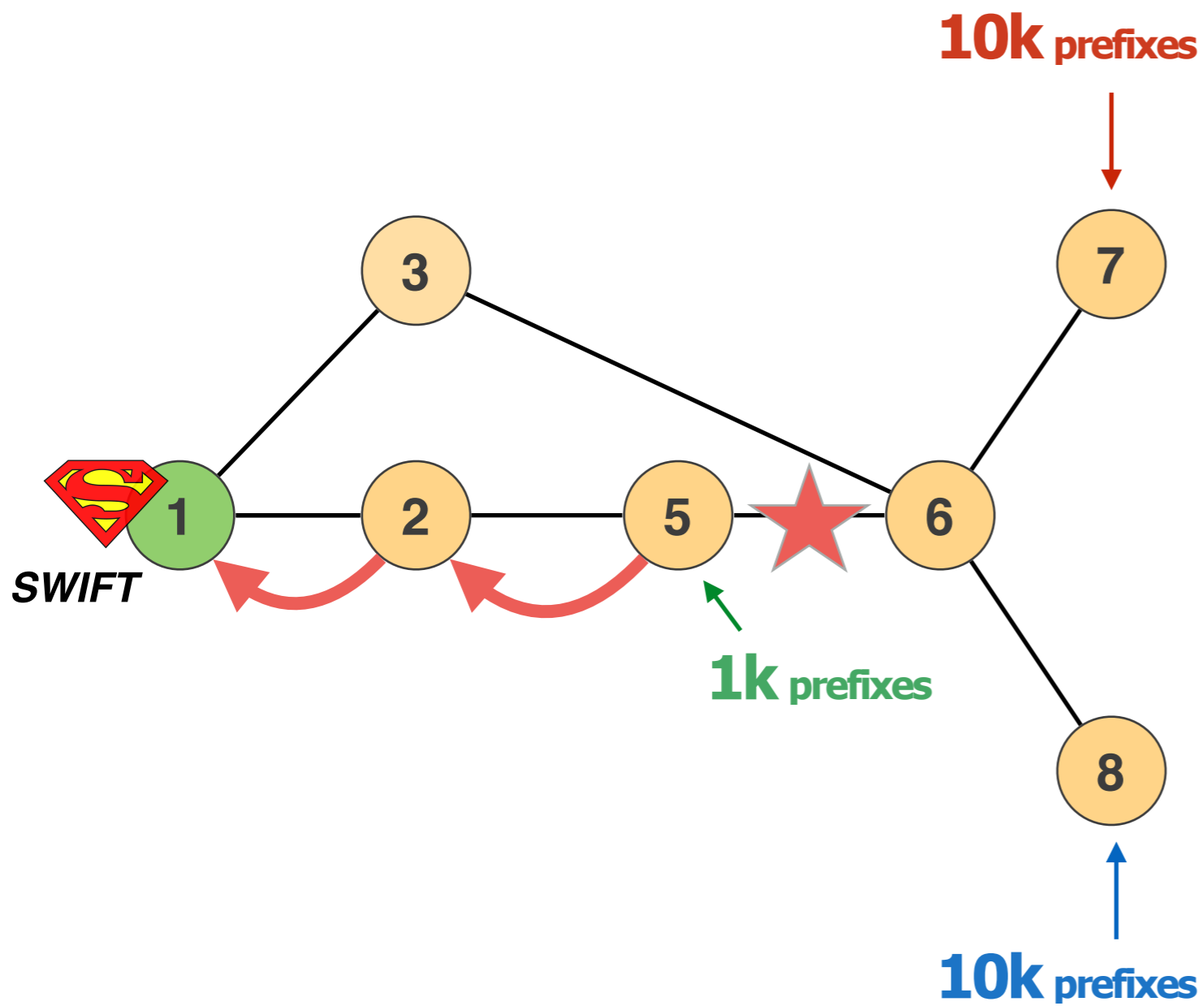
10k prefixes

AS links candidates set

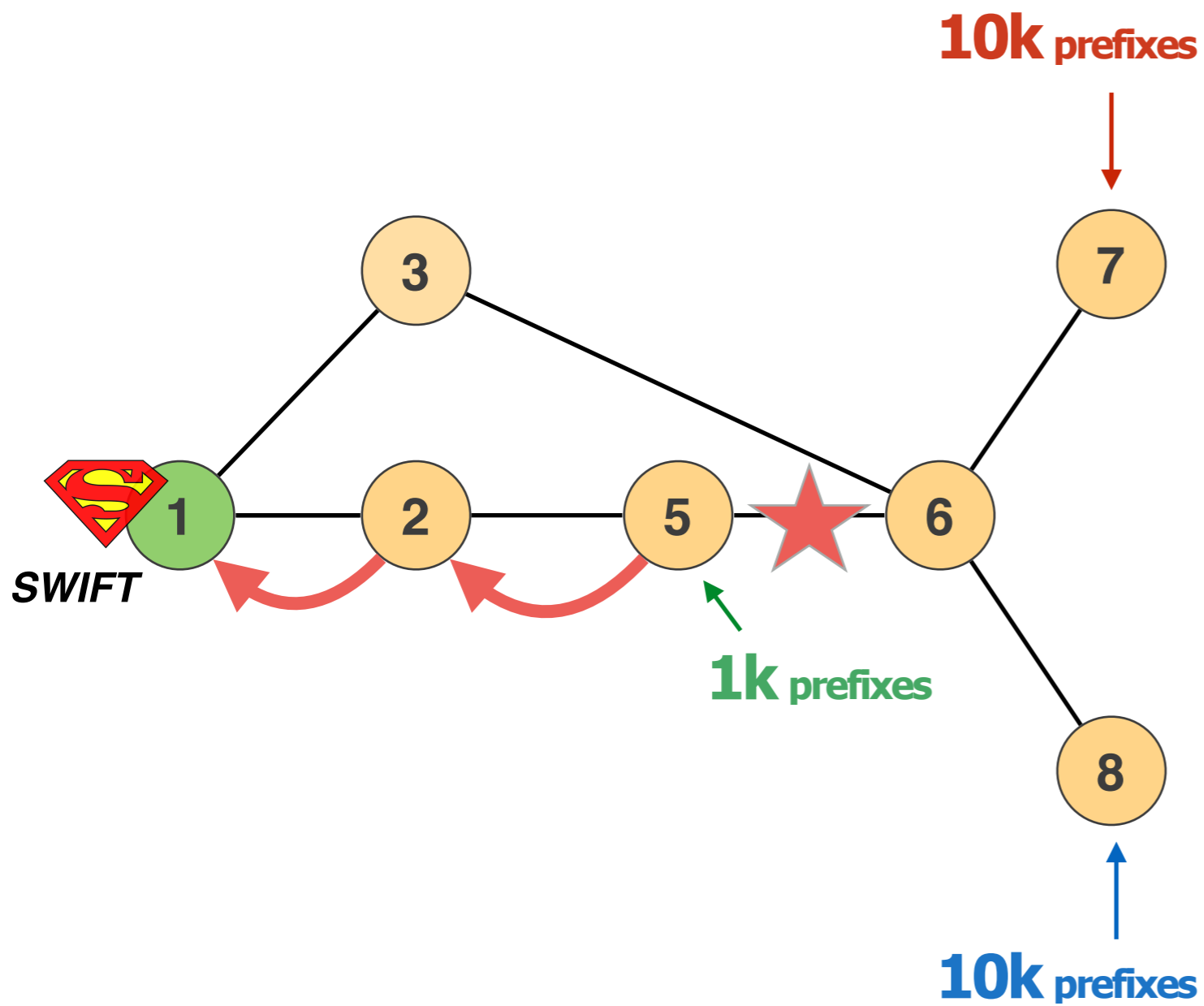
(2,5) (5,6) (6,7) (6,8)

affected prefixes

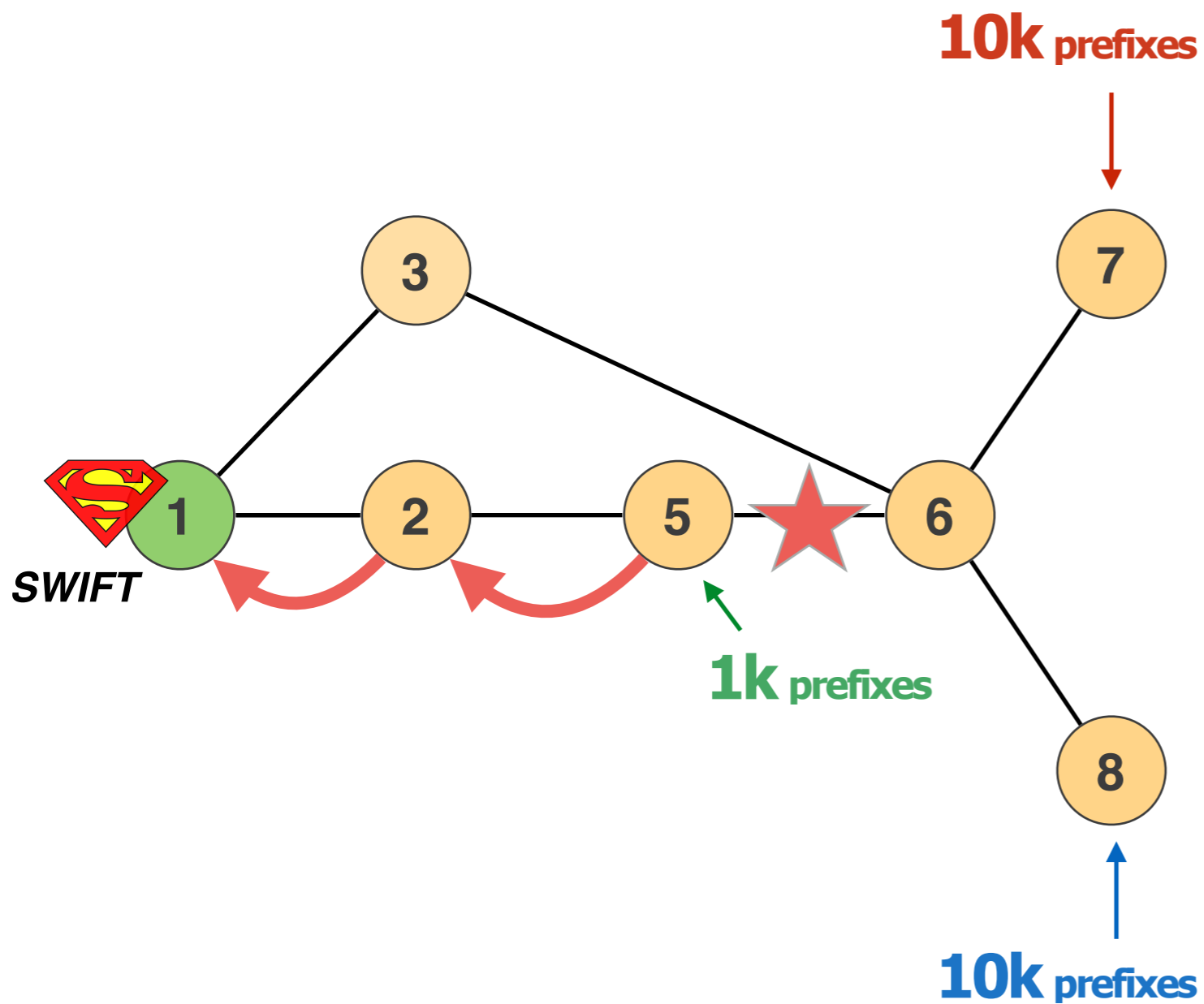
from AS7 (2,5) (5,6) (6,7) ~~(6,8)~~



	AS links	candidates set
	(2,5) (5,6) (6,7) (6,8)	(2,5) (5,6) (6,7) (6,8)
affected prefixes		
from AS7	(2,5) (5,6) (6,7) (6,8)	(2,5) (5,6) (6,7) (6,8)
from AS8	(2,5) (5,6) (6,7) (6,8)	(2,5) (5,6) (6,7) (6,8)



	candidates set
AS links	(2,5) (5,6) (6,7) (6,8)
affected prefixes	
from AS7	(2,5) (5,6) (6,7) (6,8)
from AS8	(2,5) (5,6) (6,7) (6,8)
unaffected prefixes	
from AS5	(2,5) (5,6) (6,7) (6,8)



	candidates set
AS links	(2,5) (5,6) (6,7) (6,8)
affected prefixes	
from AS7	(2,5) (5,6) (6,7) (6,8)
from AS8	(2,5) (5,6) (6,7) (6,8)
unaffected prefixes	
from AS5	(2,5) (5,6) (6,7) (6,8)

is the failed link

In practice, SWIFT uses **two metrics**
Computed on a per AS-link basis

Withdrawals share

$WS(l, t)$

**Fraction of withdrawals
crossing link l**

Path sharing

$PS(l, t)$

**Proportion of prefixes
withdrawn on link l**

In practice, SWIFT uses **two metrics**
Computed on a per AS-link basis

$$FS(l, t) = \boxed{\begin{array}{c} \text{Withdrawals share} \\ WS(l, t) \end{array}} \times \boxed{\begin{array}{c} \text{Path sharing} \\ PS(l, t) \end{array}}$$

Fit Score for link l Fraction of withdrawals crossing link l Proportion of prefixes withdrawn on link l

Theorem

Under perfect conditions, SWIFT **always** returns a set of links including the failed link

Theorem

Under perfect conditions, SWIFT **always** returns a set of links including the failed link

Challenge #1

To be fast, we can't wait for all the withdrawals

SWIFT runs the link inference **early** during the burst in order to predict future withdrawals

Theorem

Under perfect conditions, SWIFT **always** returns a set of links including the failed link

Challenge #1

To be fast, we can't wait for all the withdrawals

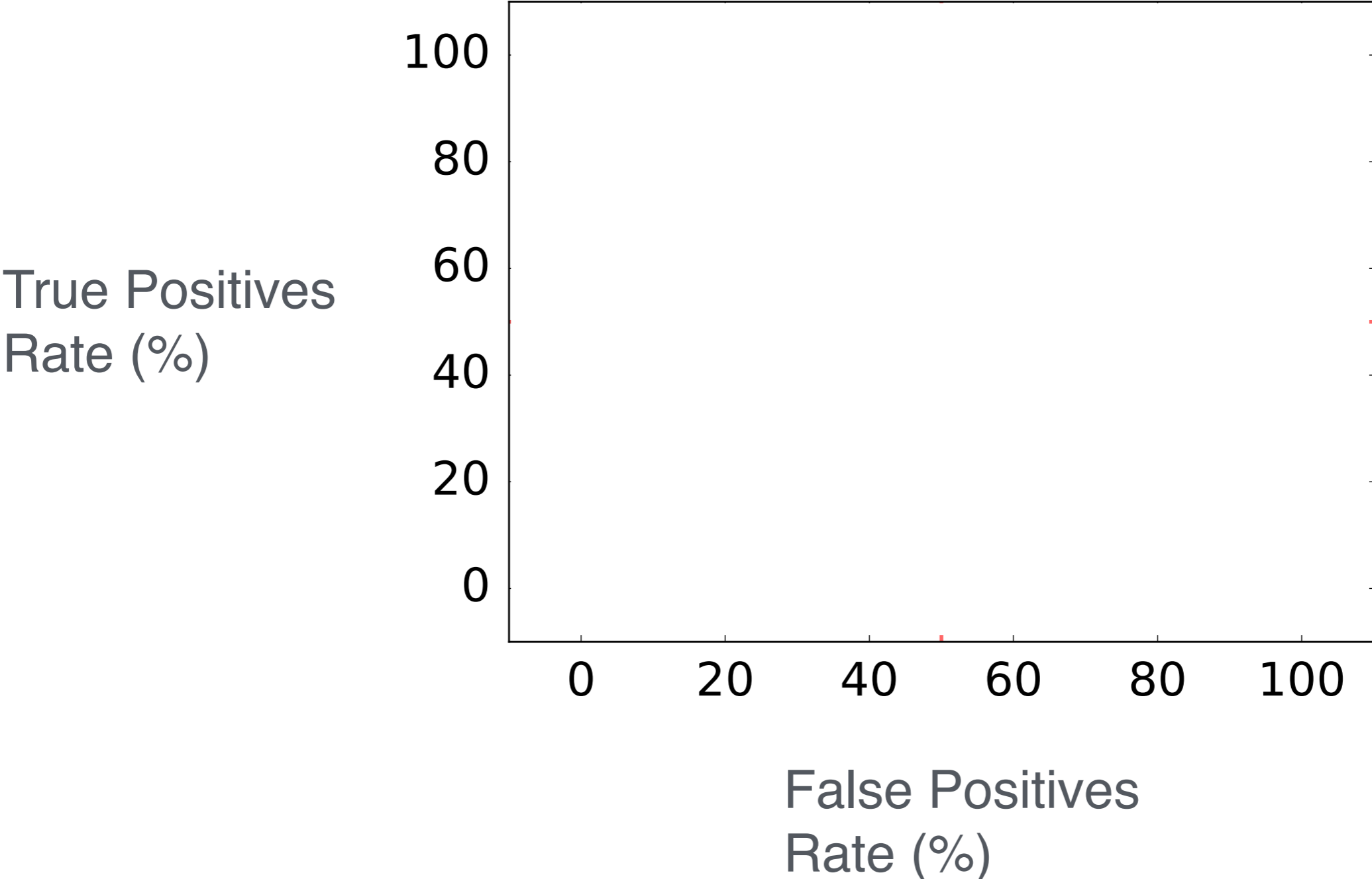
SWIFT runs the link inference **early** during the burst in order to predict future withdrawals

Challenge #2

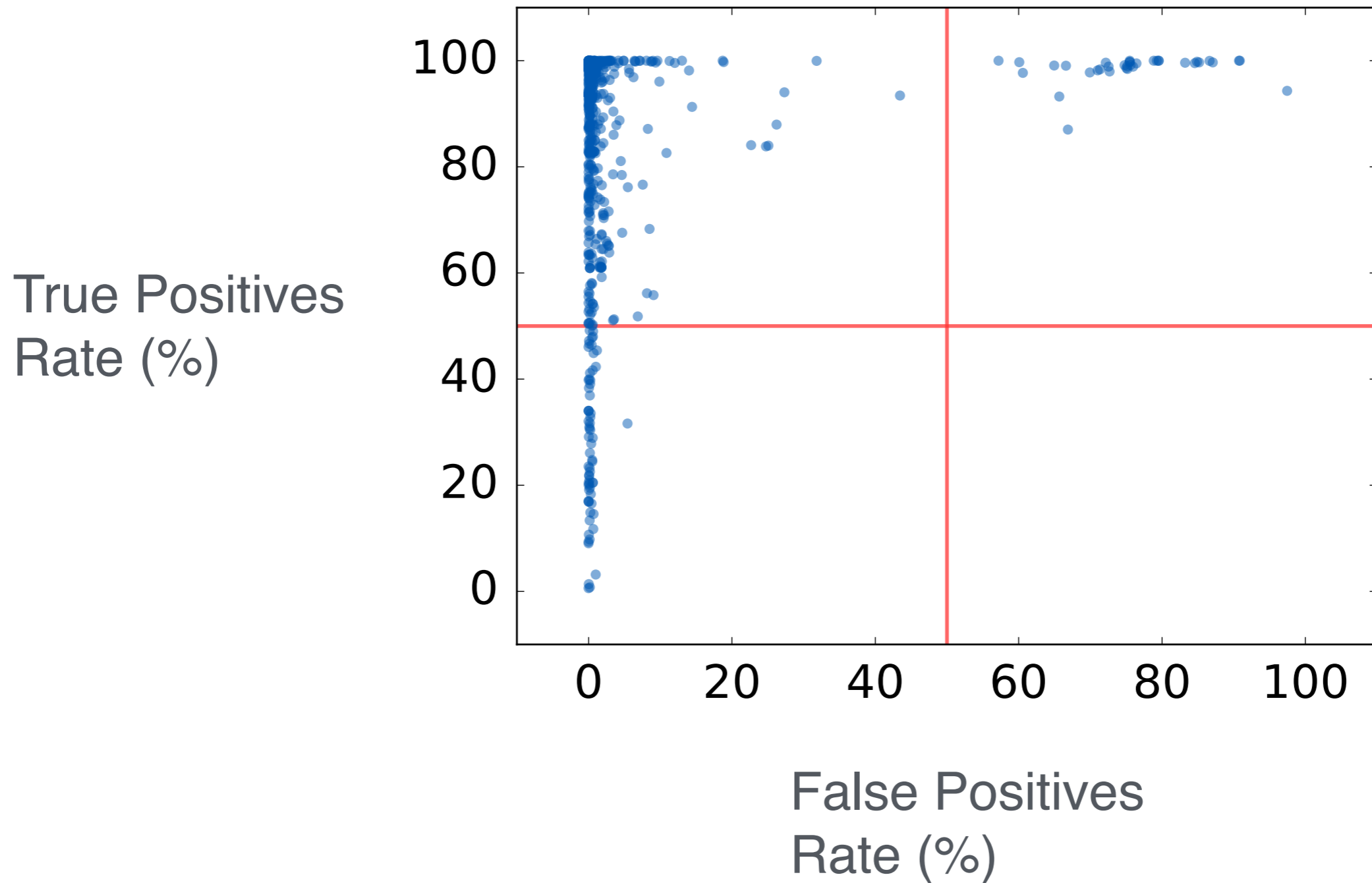
An outage can affect multiple AS links

*SWIFT link inference algorithm returns a **set of links**, all with a high fit score*

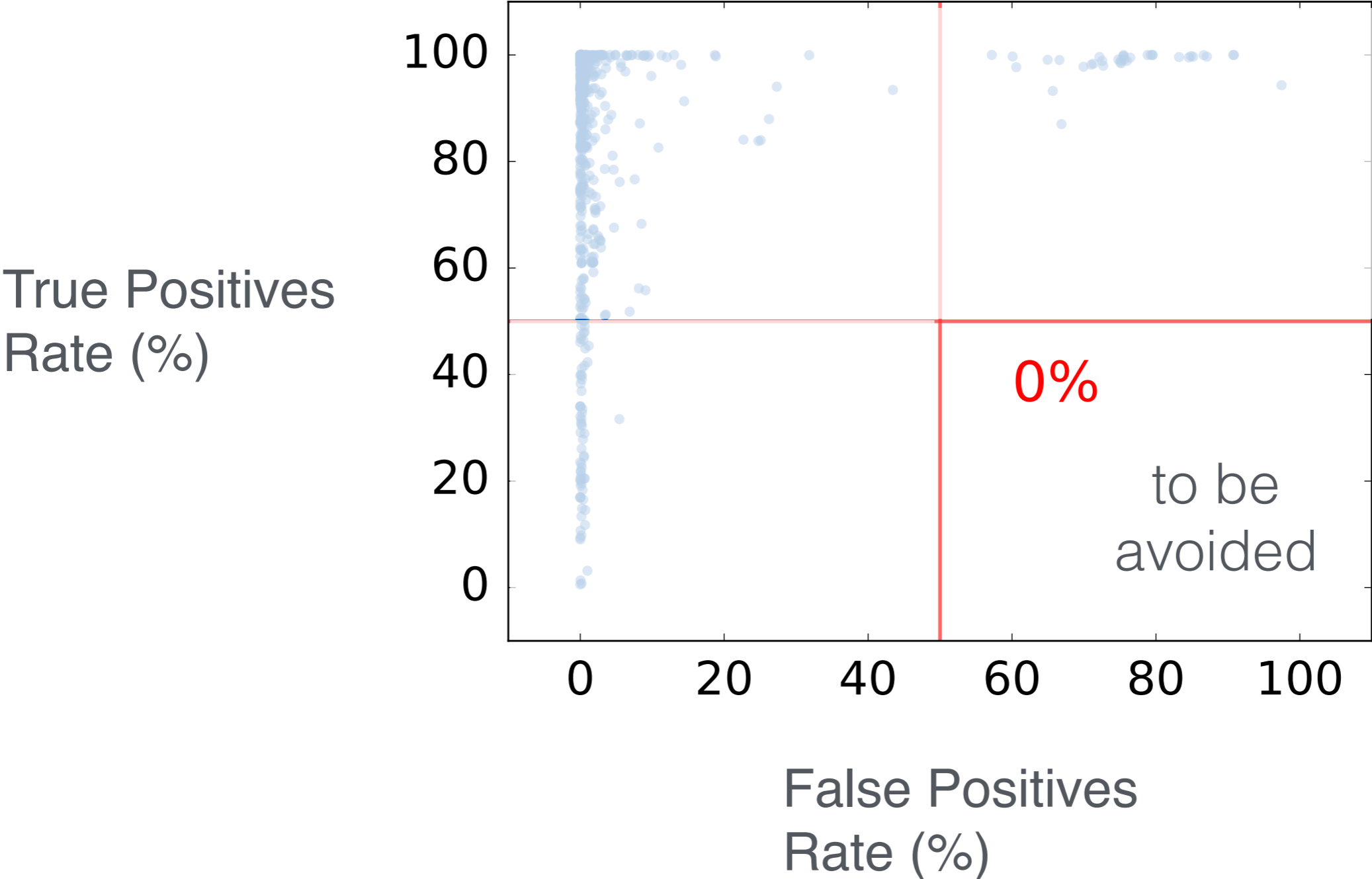
SWIFT link failure inference **is accurate**



SWIFT link failure inference **is accurate**



SWIFT link failure inference **is accurate**



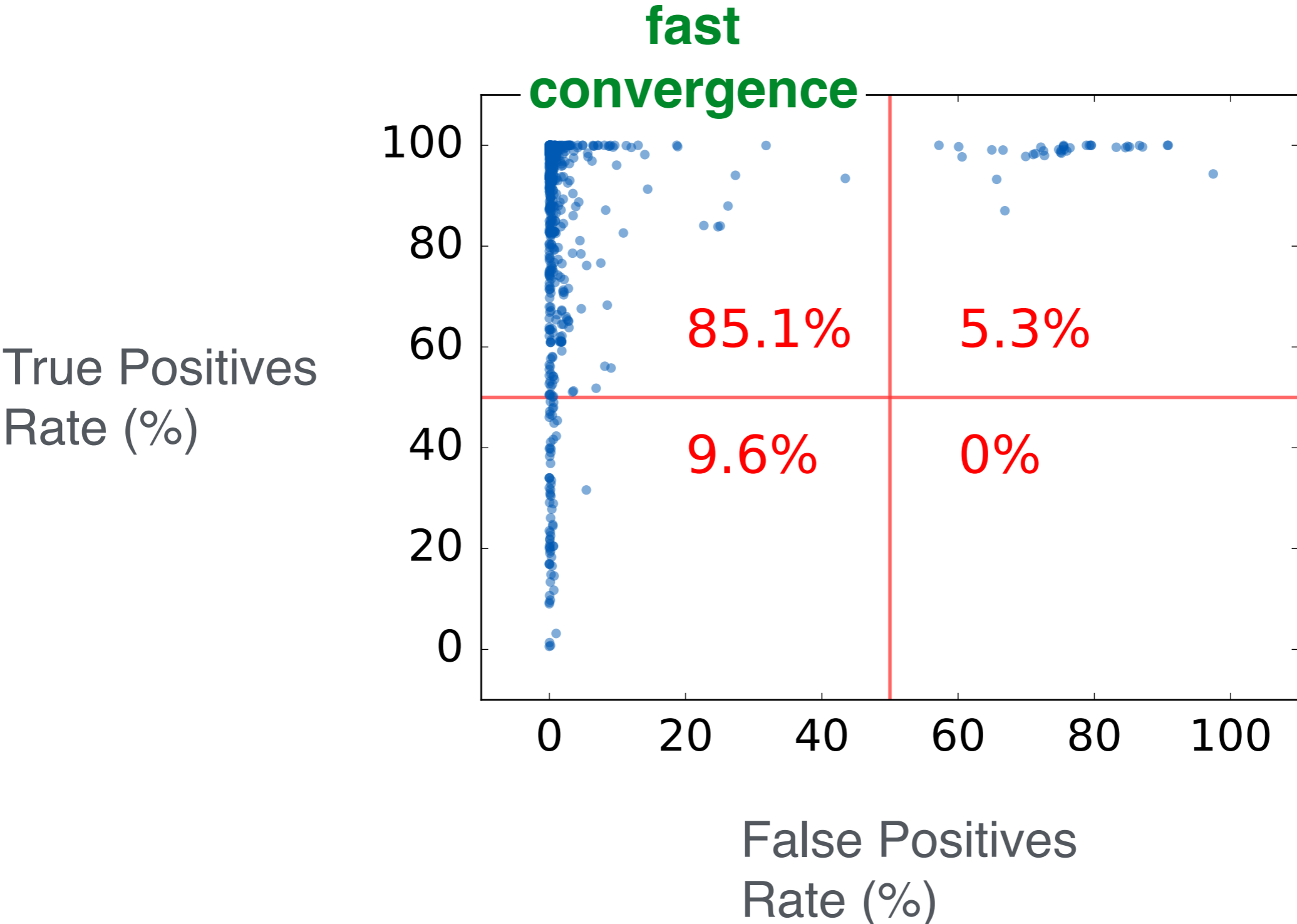
SWIFT link failure inference **is accurate**



SWIFT link failure inference **is accurate**



SWIFT link failure inference **is accurate**



SWIFT is **much faster** than BGP

	BGP	SWIFT
learning time (median)	13 sec	

SWIFT is **much faster** than BGP

	BGP	SWIFT
learning time (median)	13 sec	2 sec

SWIFT: Predictive Fast Reroute

1. **SWIFT** reduces the learning time of a withdrawal from **13s to 2s** using a control-plane prediction algorithm
2. **SWIFT** fast reroutes the affected traffic towards unaffected paths in **40ms** using a hierarchical forwarding table
3. **SWIFT** is **deployable** on existing routers

SWIFT: Predictive Fast Reroute

1. ***SWIFT*** reduces the learning time of a withdrawal from **13s to 2s** using a control-plane prediction algorithm
2. ***SWIFT*** fast reroutes the affected traffic towards unaffected paths in **40ms** using a hierarchical forwarding table
3. ***SWIFT*** is **deployable** on existing routers

Fast rerouting the traffic upon remote outages
has several challenges

Fast rerouting the traffic upon remote outages has several challenges

A remote outage can affect any set of prefixes

Fast rerouting the traffic upon remote outages has several challenges

A remote outage can affect any set of prefixes

A remote outage can affect backup paths

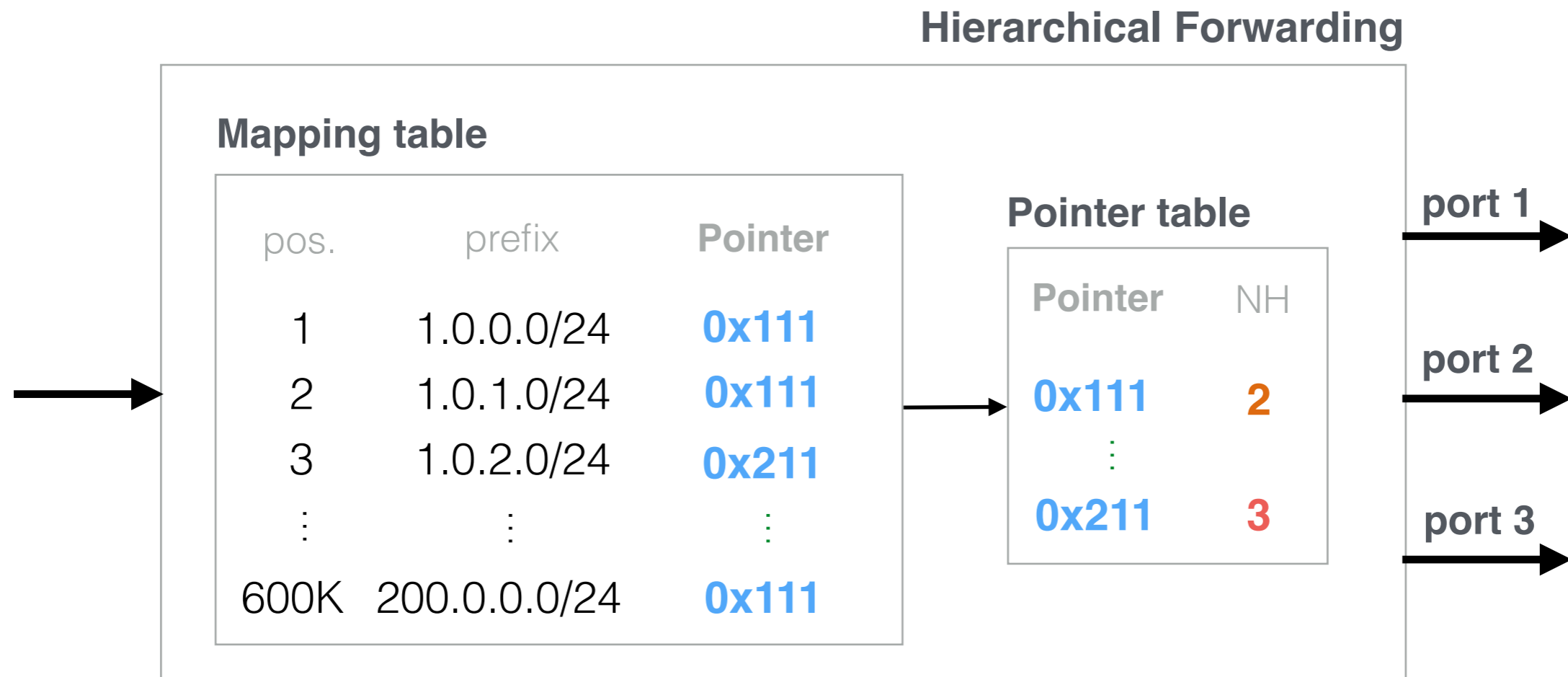
Fast rerouting the traffic upon remote outages has several challenges

A remote outage can affect any set of prefixes

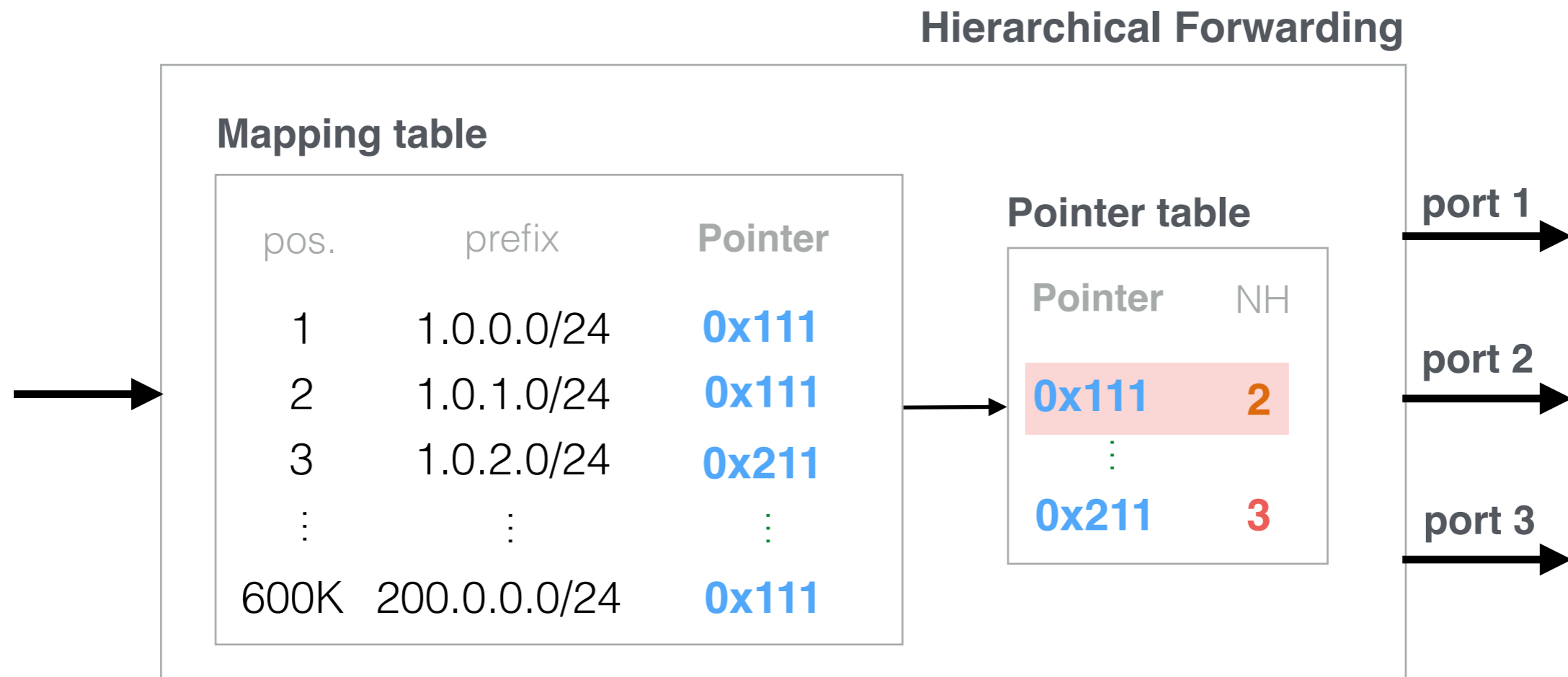
A remote outage can affect backup paths

and we want to be fast

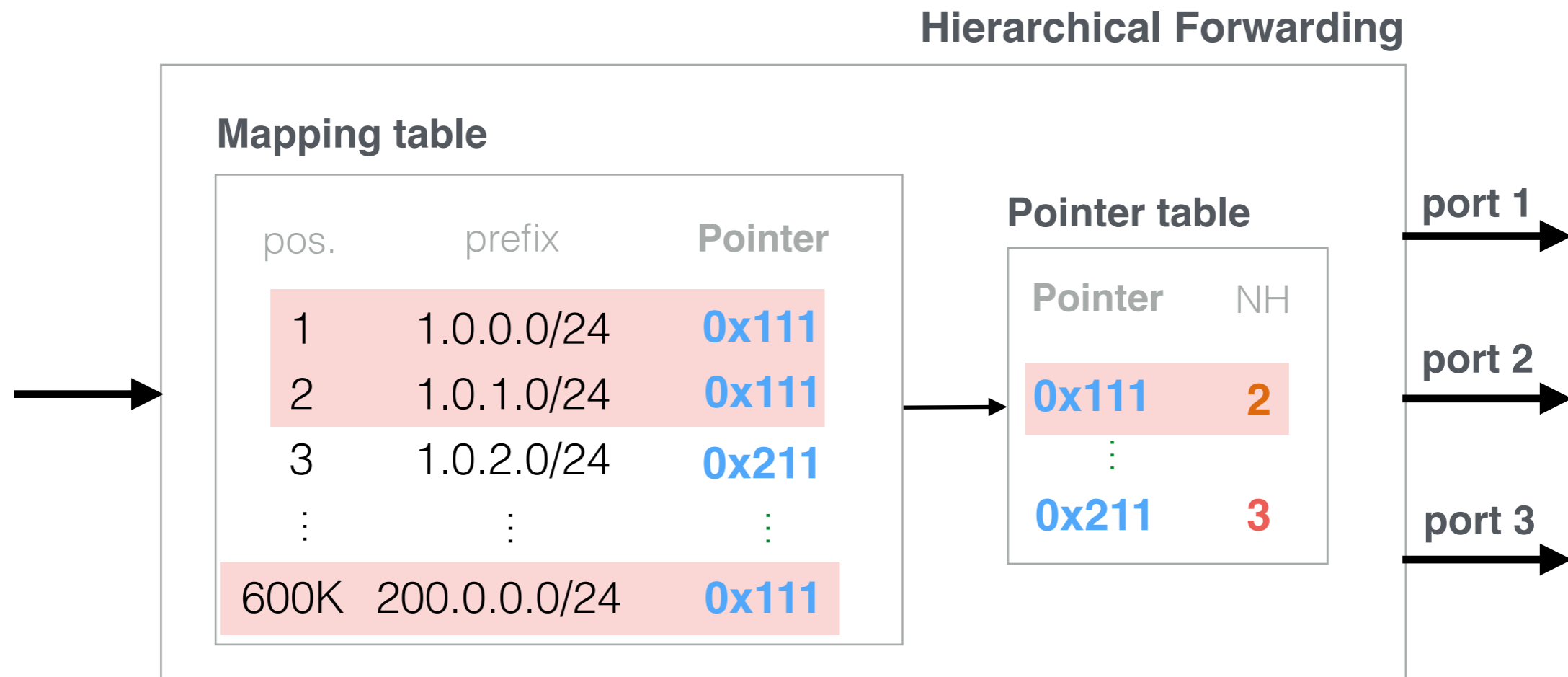
To solve these challenges,
SWIFT relies on a **hierarchical forwarding table**



To solve these challenges,
SWIFT relies on a **hierarchical forwarding table**

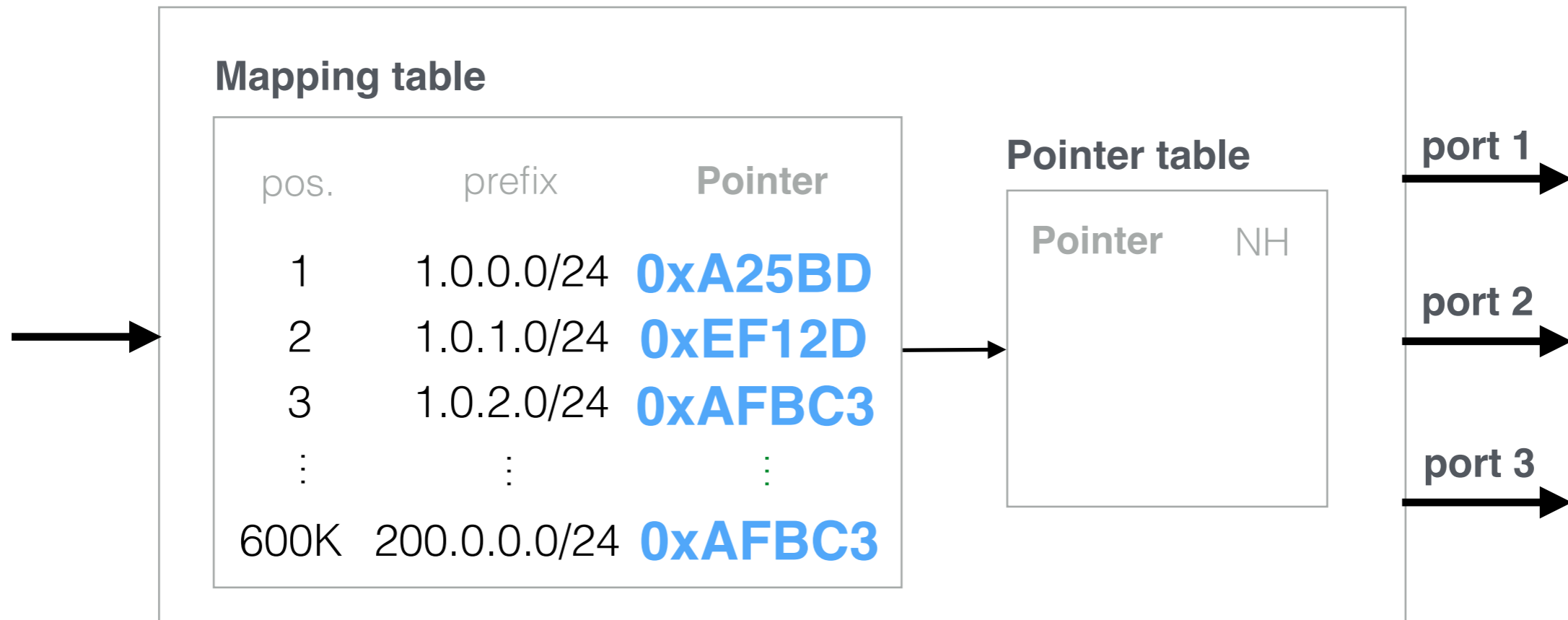


To solve these challenges,
SWIFT relies on a **hierarchical forwarding table**



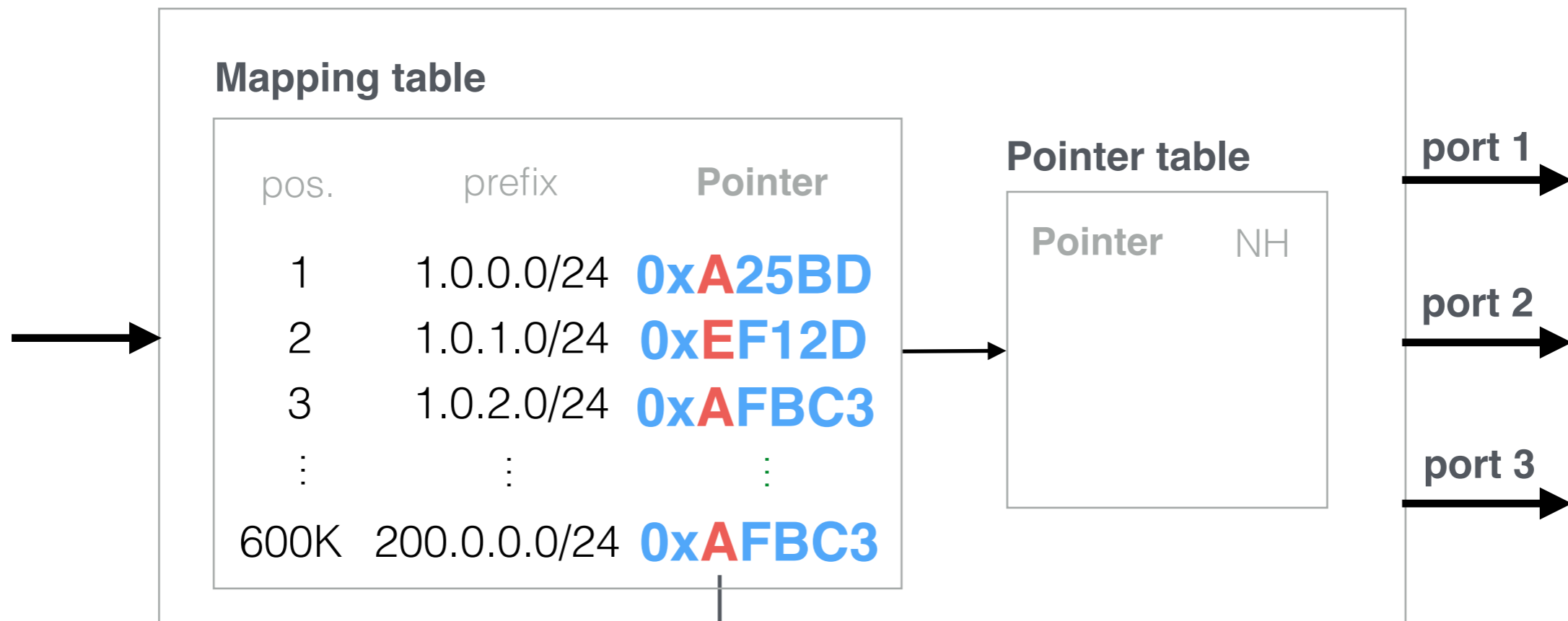
SWIFT uses custom pointers

Hierarchical Forwarding



SWIFT uses custom pointers

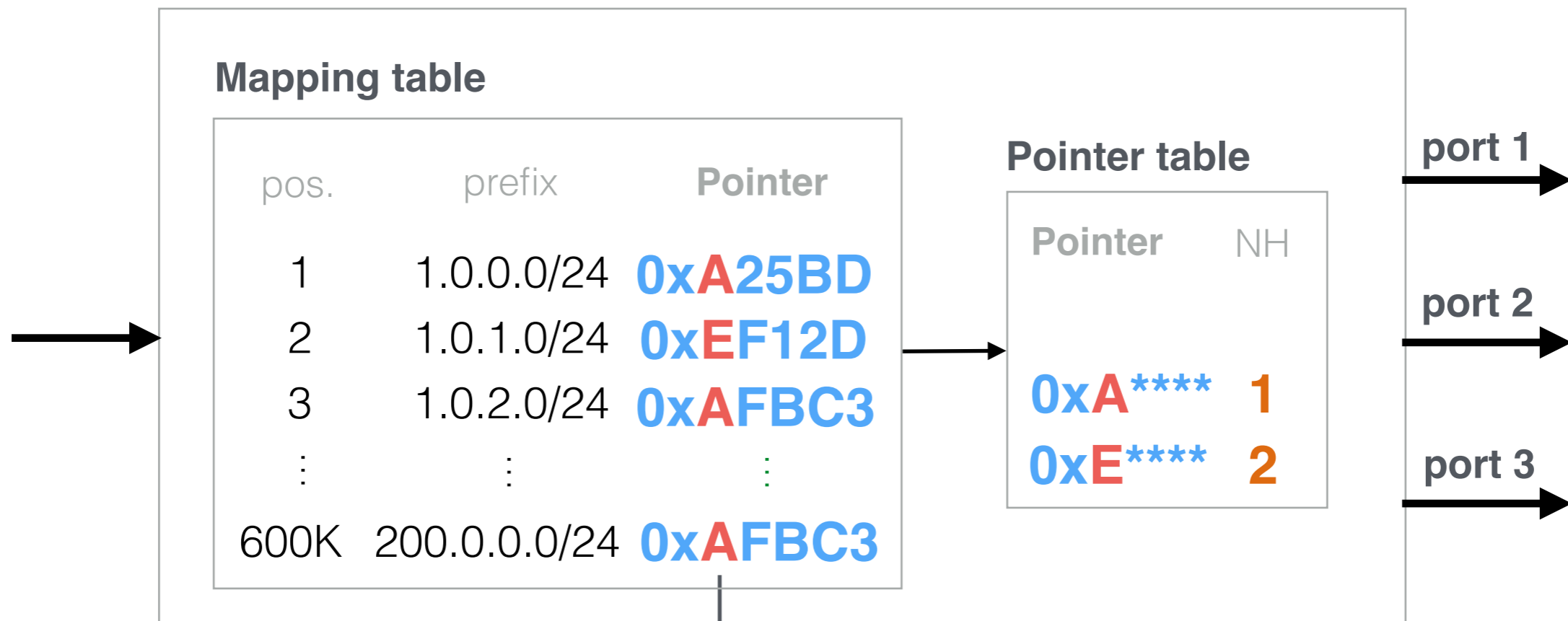
Hierarchical Forwarding



A indicates a primary next-hop

SWIFT uses custom pointers

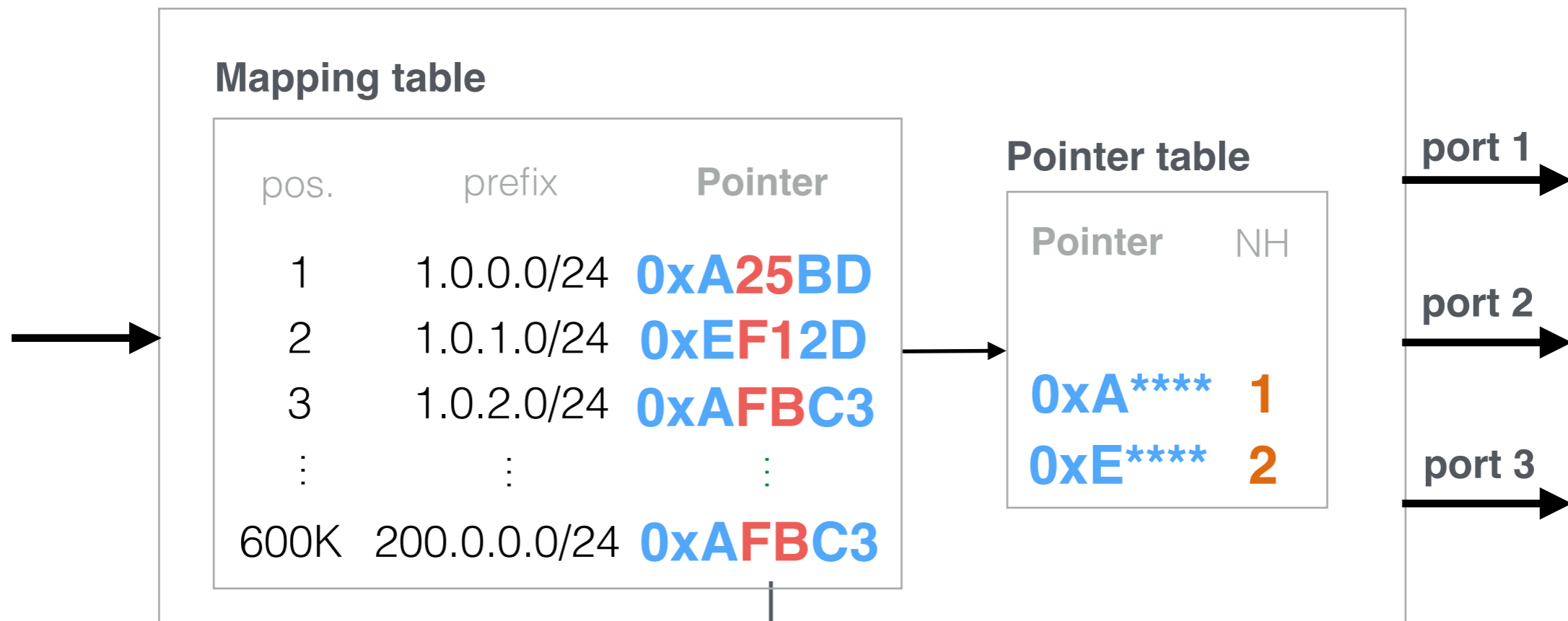
Hierarchical Forwarding



A indicates a primary next-hop

SWIFT uses custom pointers

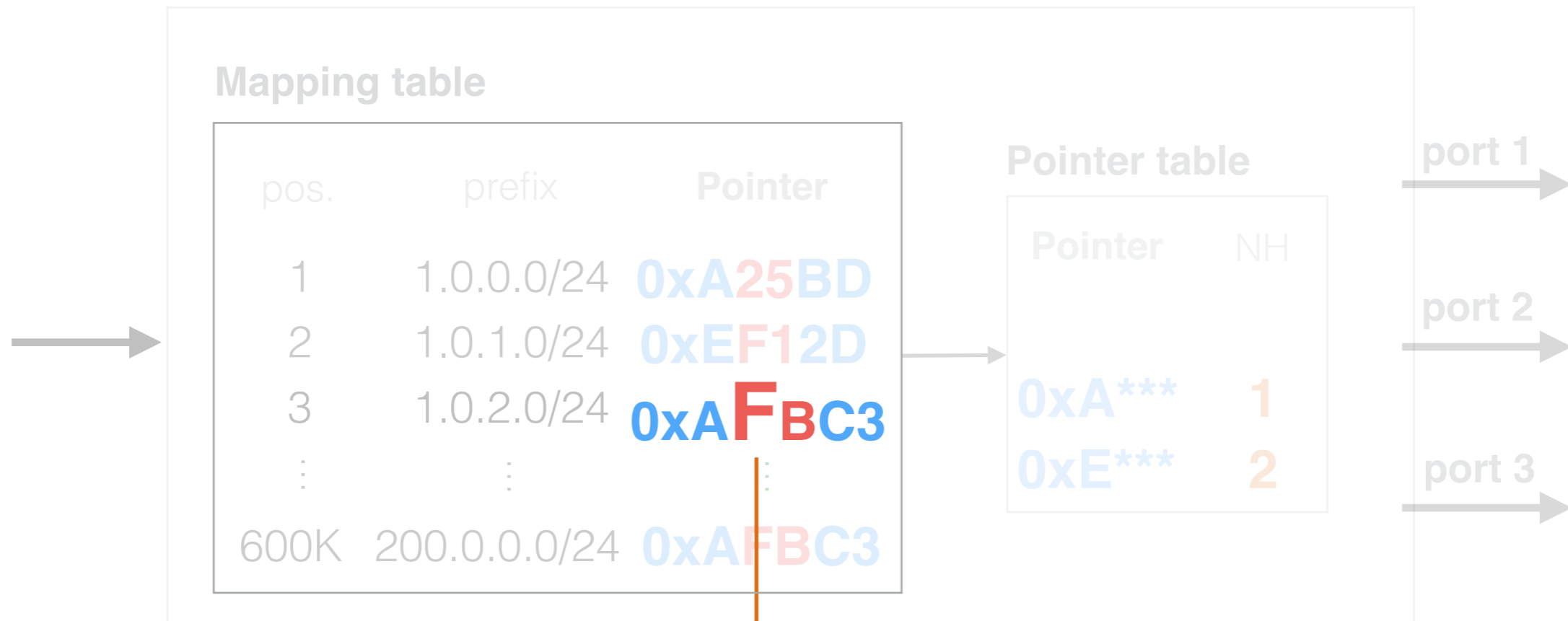
Hierarchical Forwarding



indicates the
AS path used

SWIFT uses custom pointers

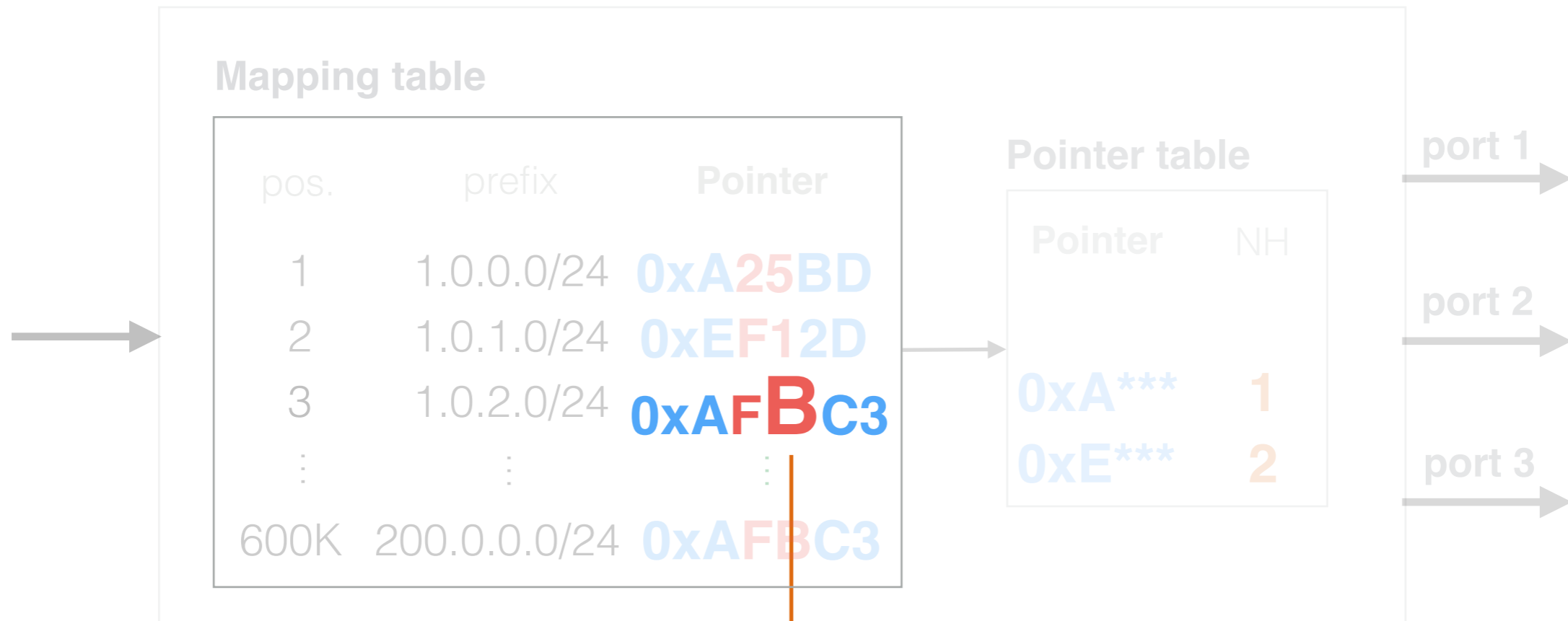
Hierarchical Forwarding



**indicates this prefix
is traversing link (5,6)**

SWIFT uses custom pointers

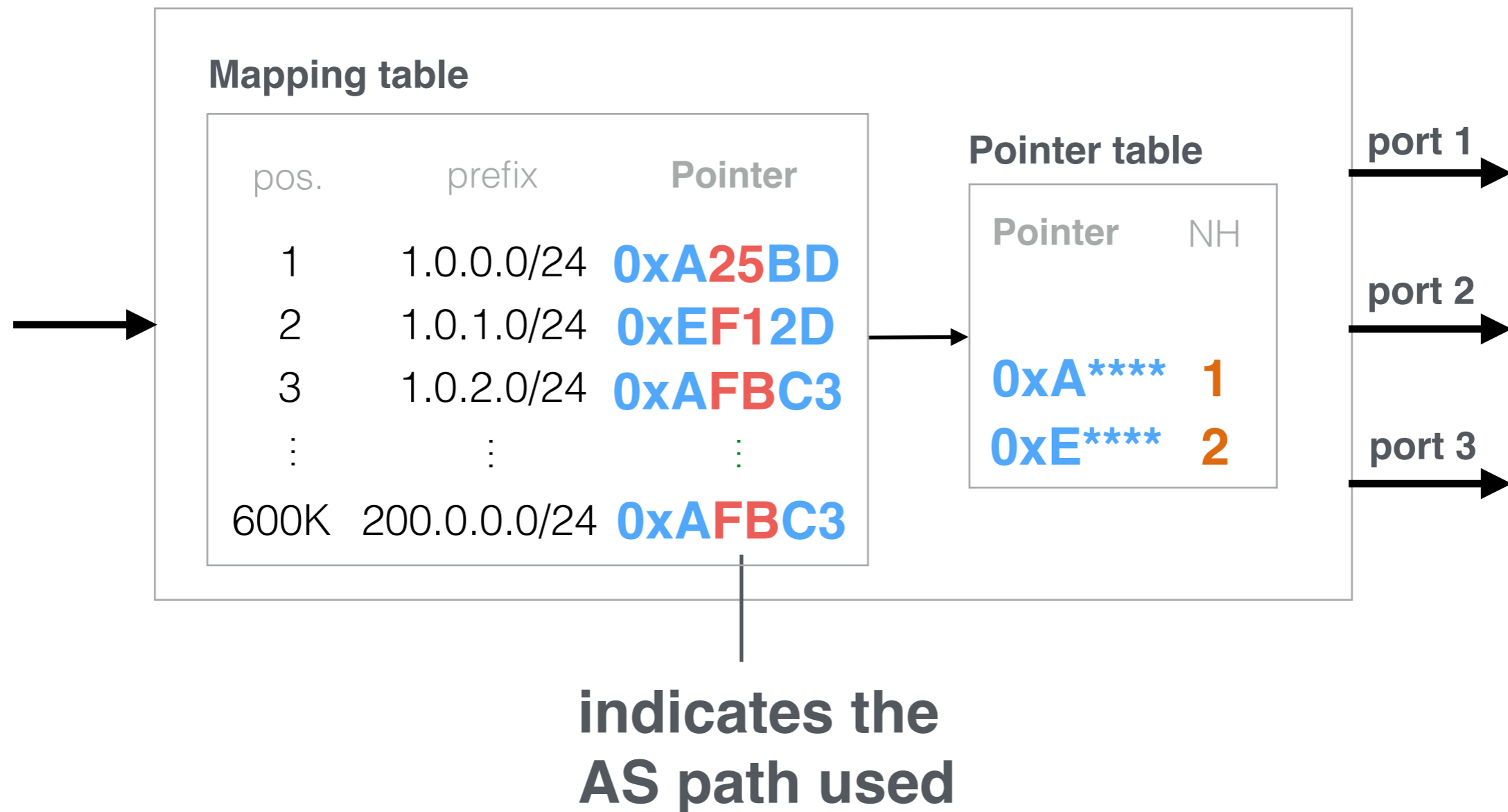
Hierarchical Forwarding



**indicates this prefix
is traversing link (6,7)**

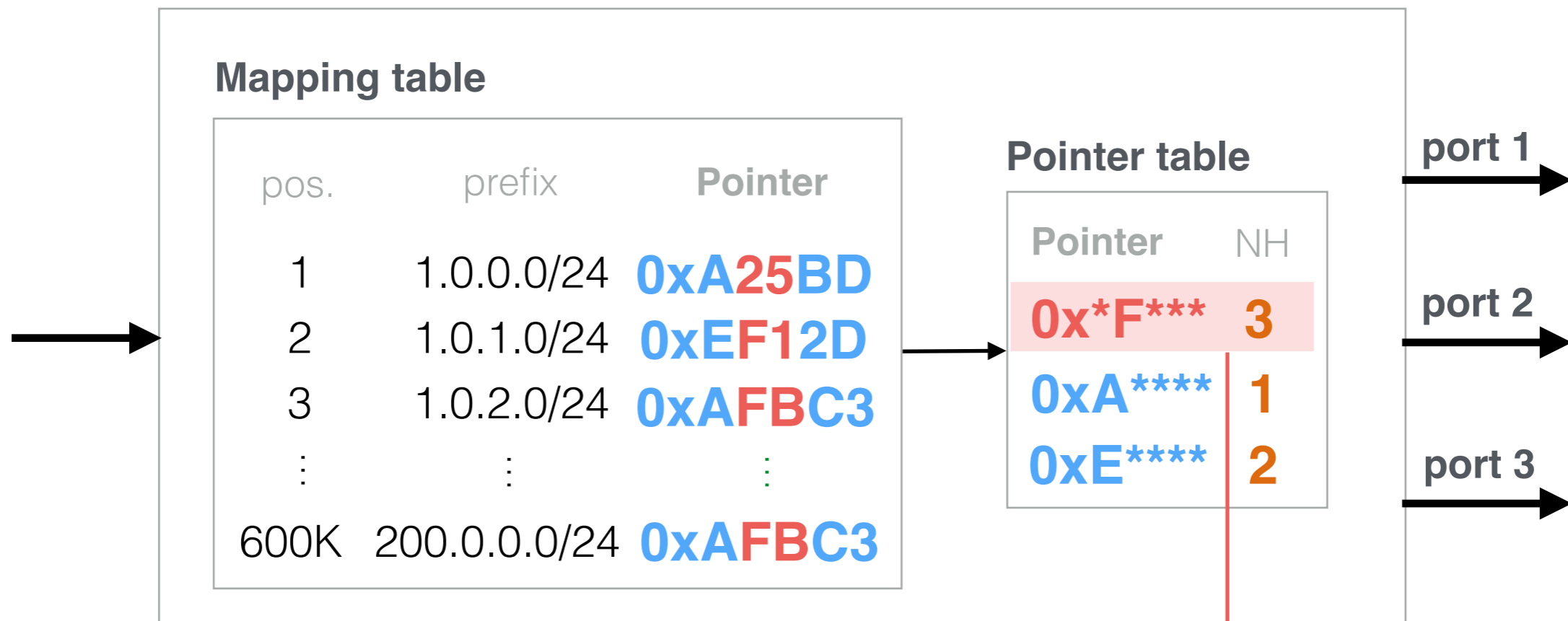
SWIFT uses custom pointers

Hierarchical Forwarding



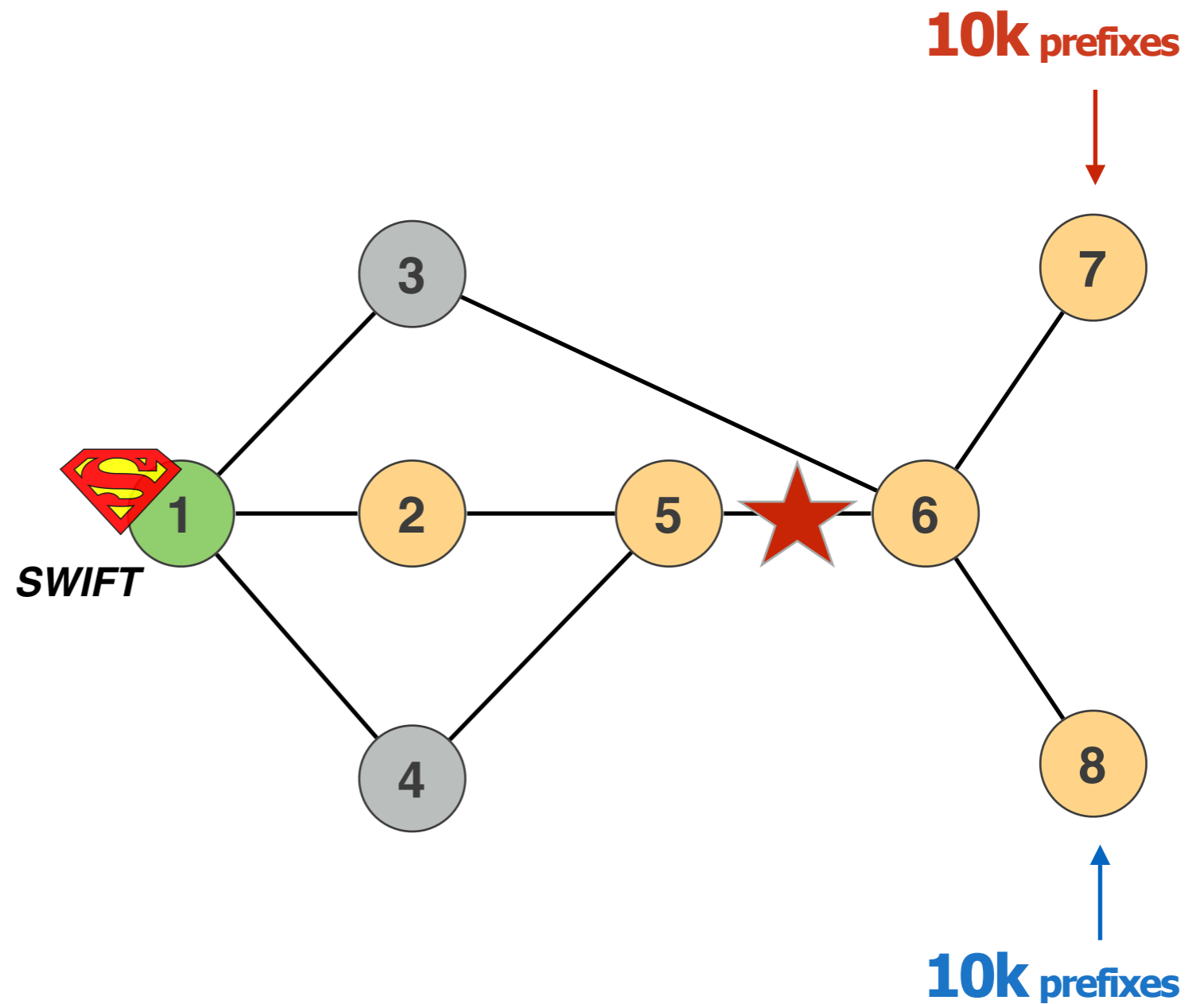
SWIFT uses custom pointers

Hierarchical Forwarding

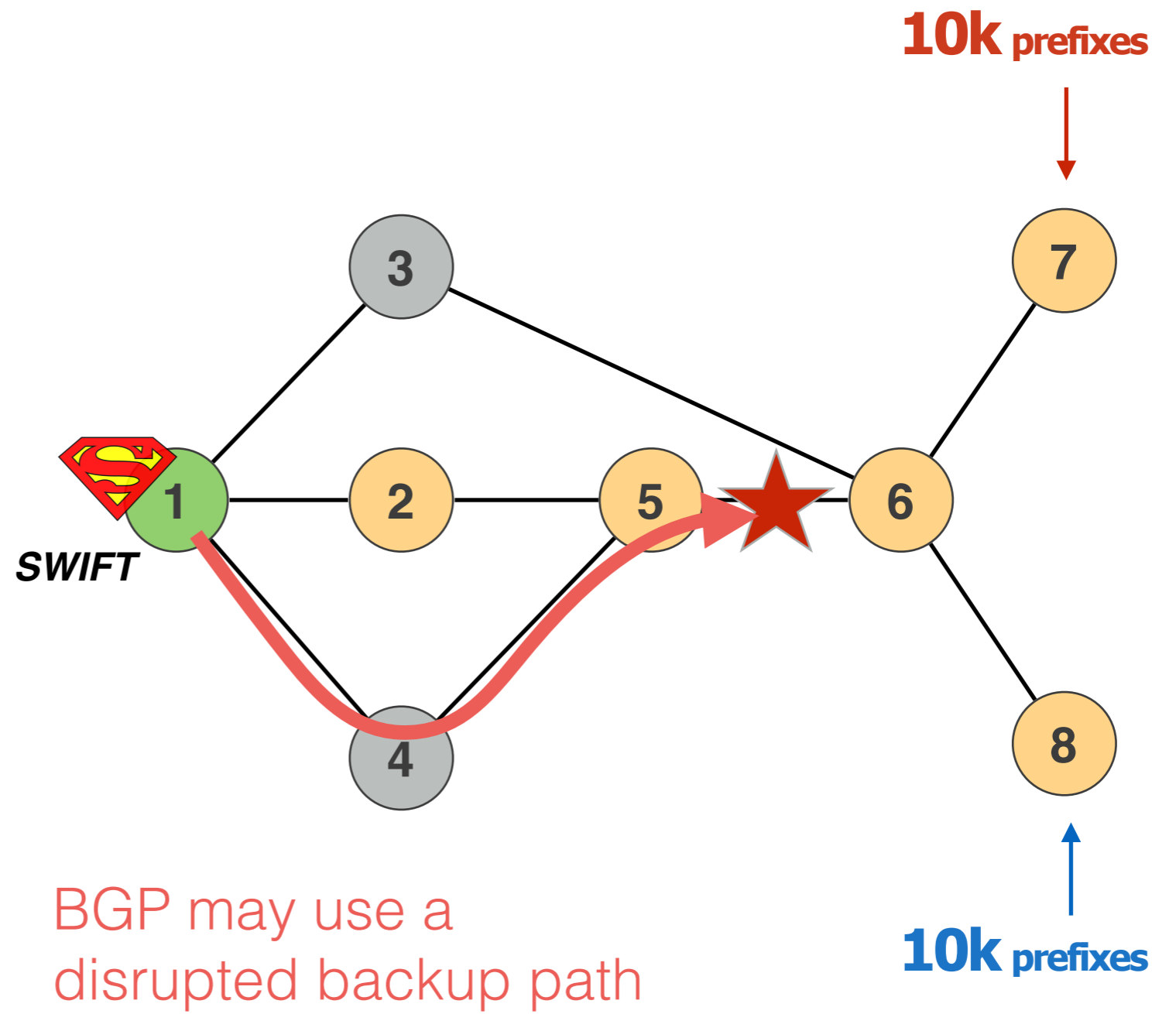


fast reroutes the traffic traversing the link (5,6)

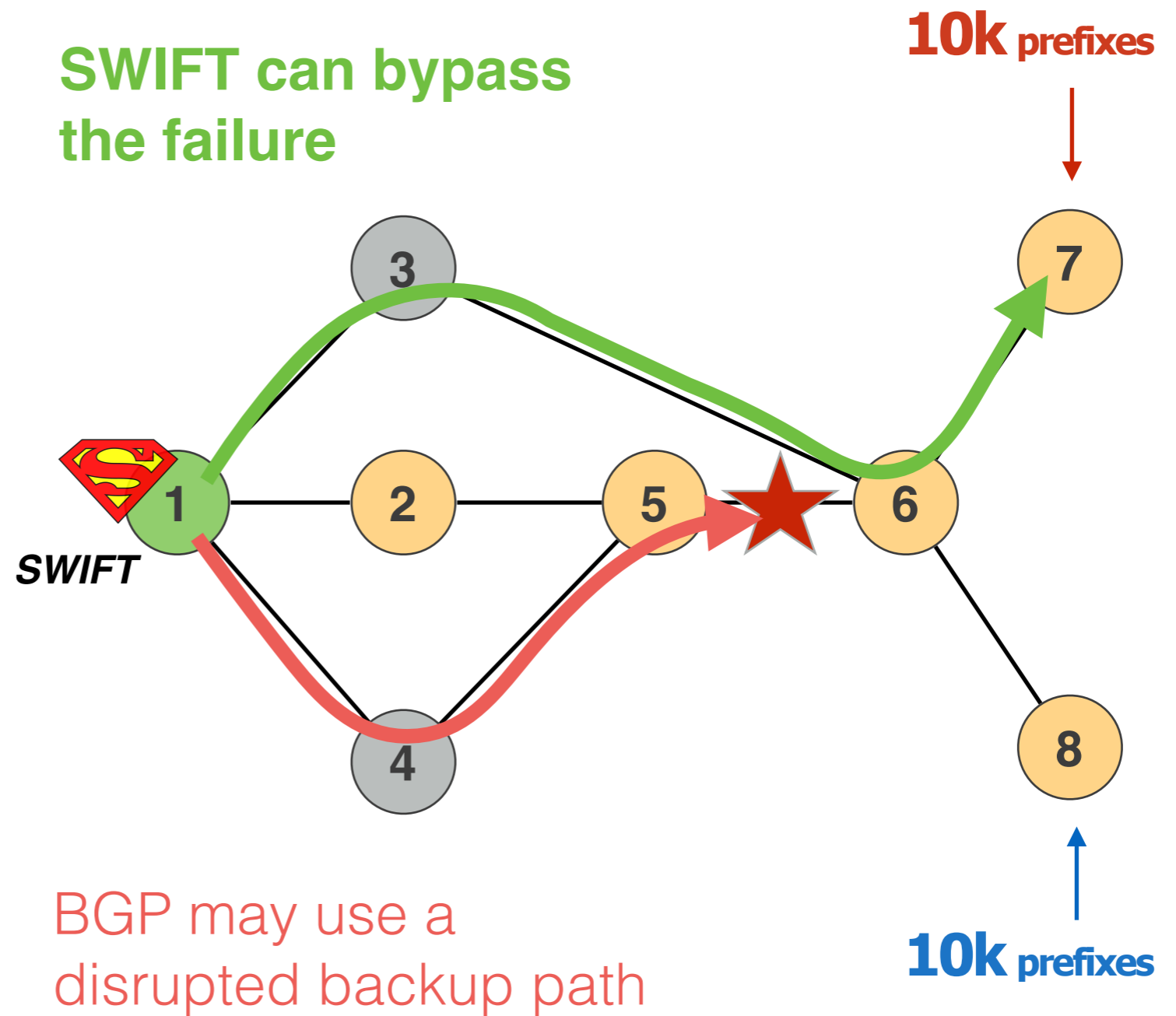
SWIFT reroutes the affected traffic to **unaffected backup paths**



SWIFT reroutes the affected traffic to **unaffected backup paths**

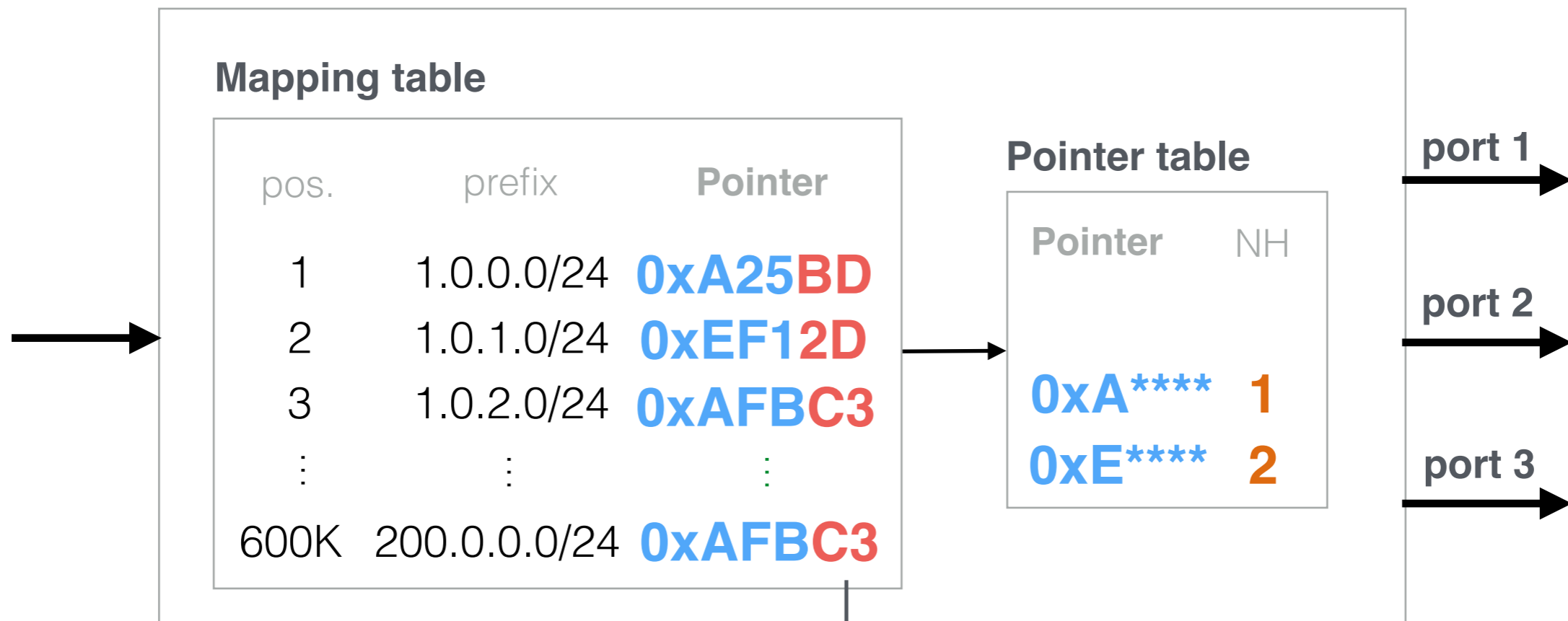


SWIFT reroutes the affected traffic to **unaffected backup paths**



SWIFT encodes the next-hops in the pointers too

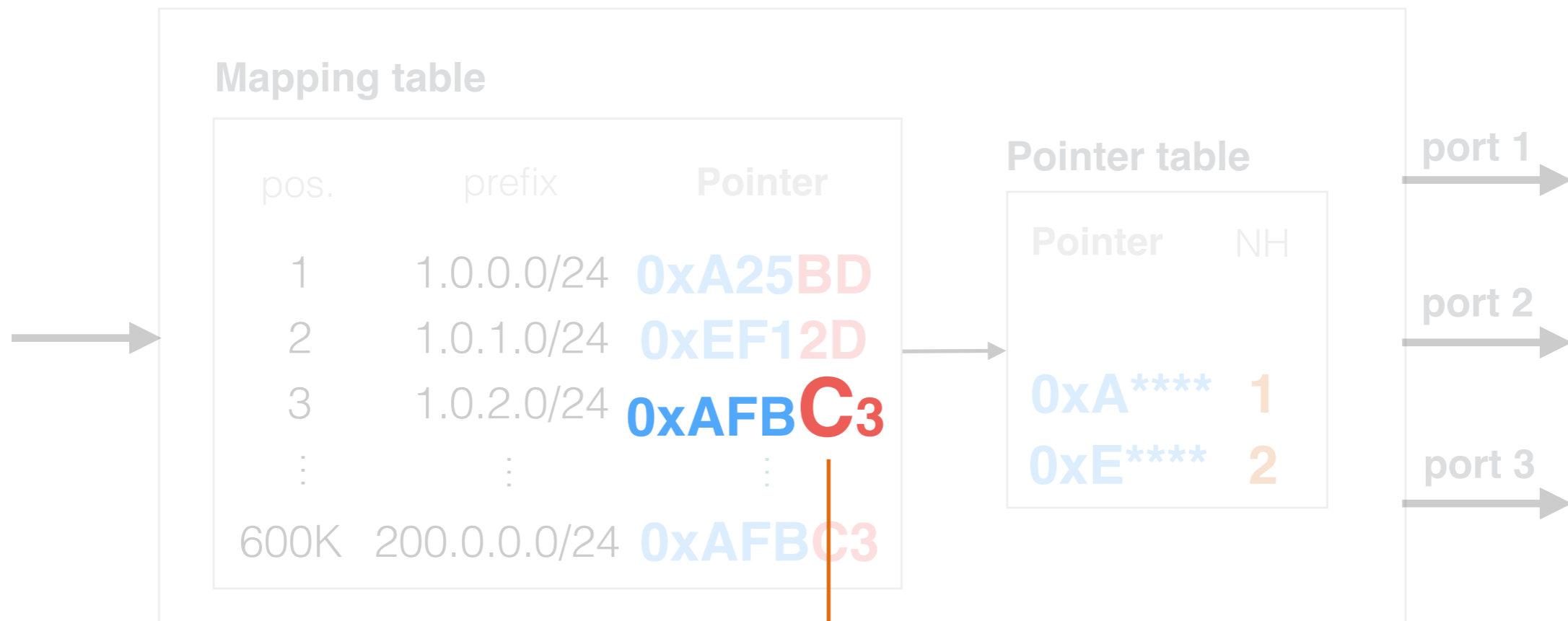
Hierarchical Forwarding



indicates the backup next-hops

SWIFT encodes the next-hops in the pointers too

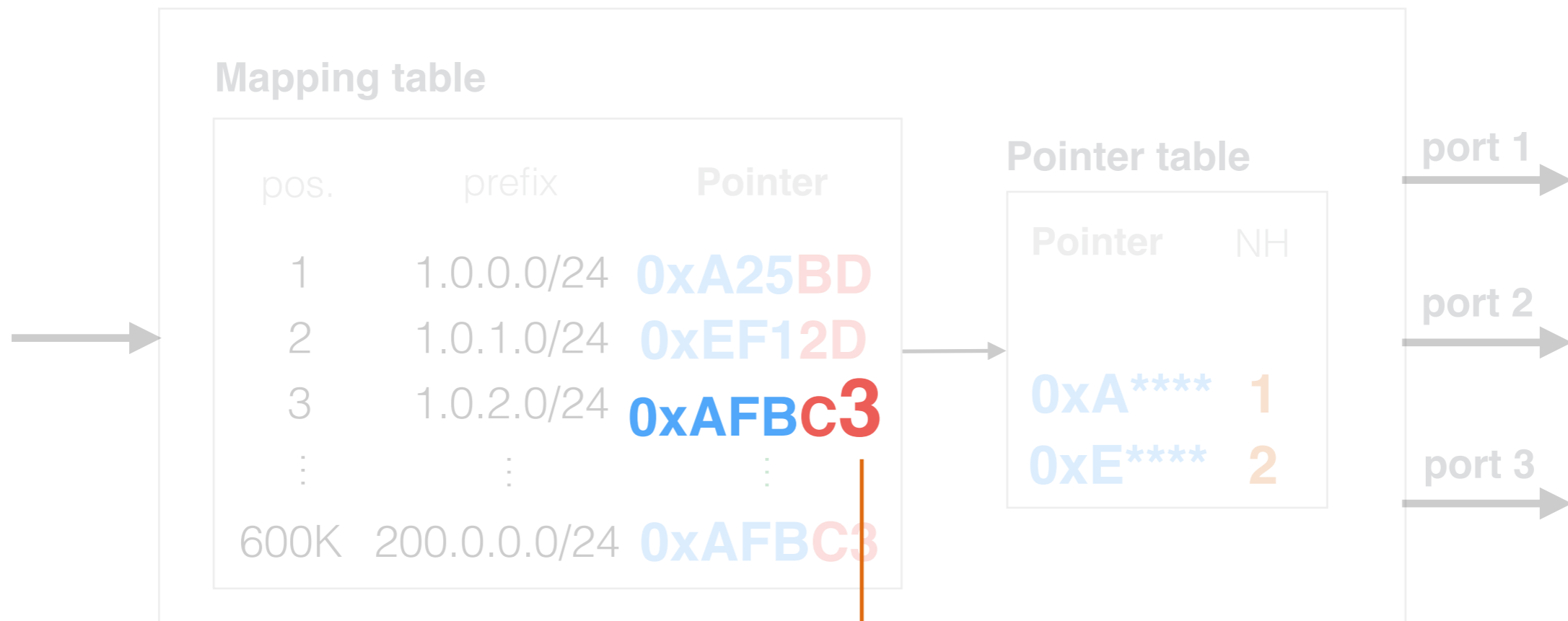
Hierarchical Forwarding



indicates the backup next-hop to use if (5,6) fails

SWIFT encodes the next-hops in the pointers too

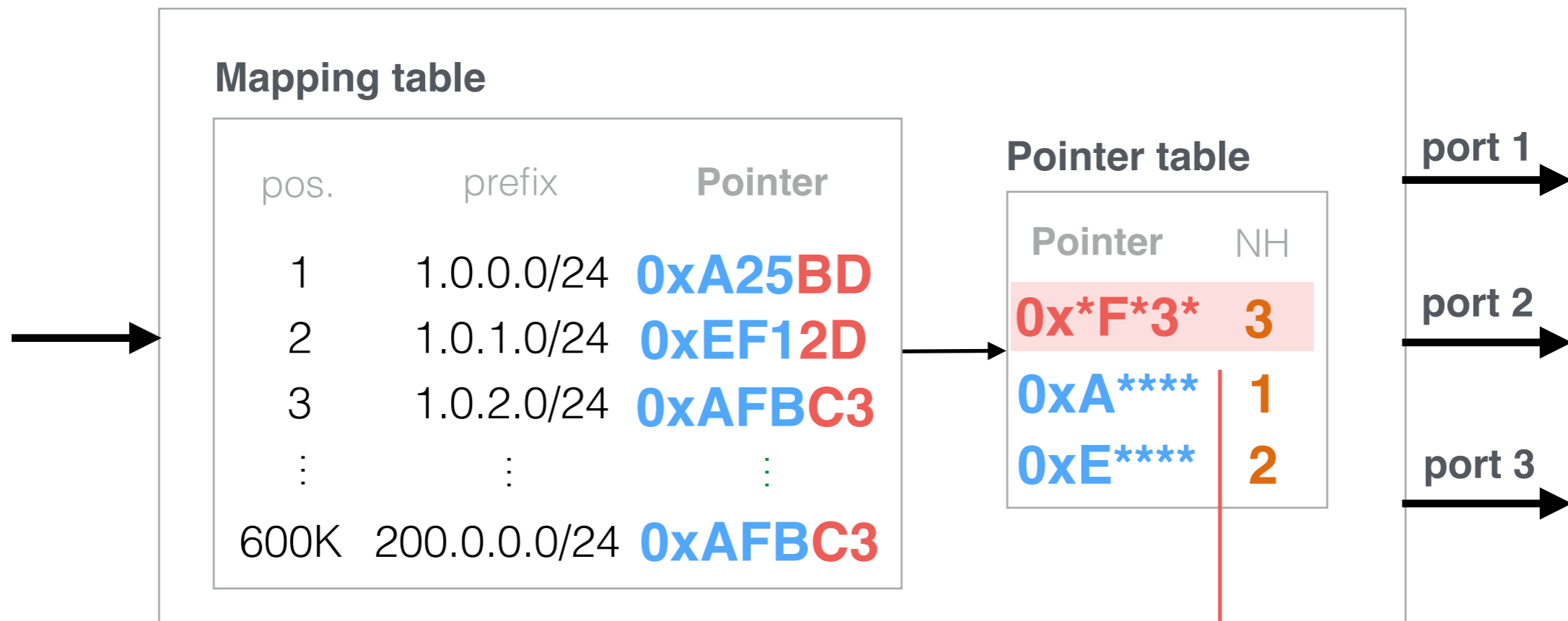
Hierarchical Forwarding



indicates the backup next-hop to use if (6,7) fails

SWIFT encodes the next-hops in the pointers too

Hierarchical Forwarding

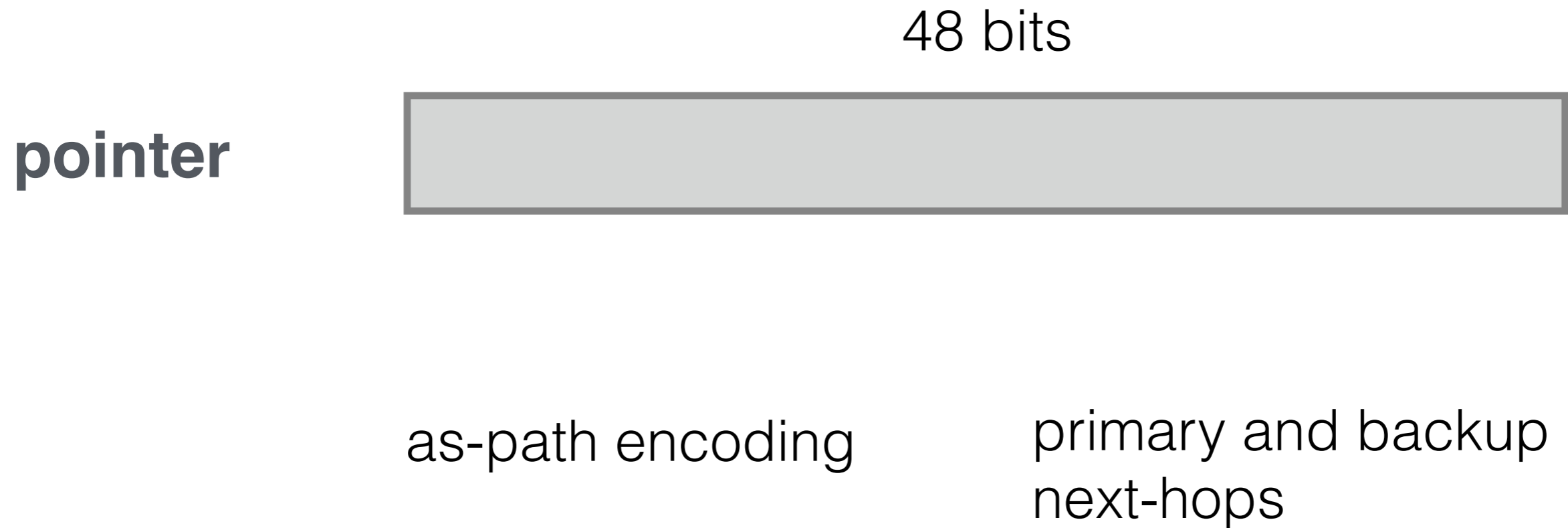


fast reroutes the traffic towards a valid backup path

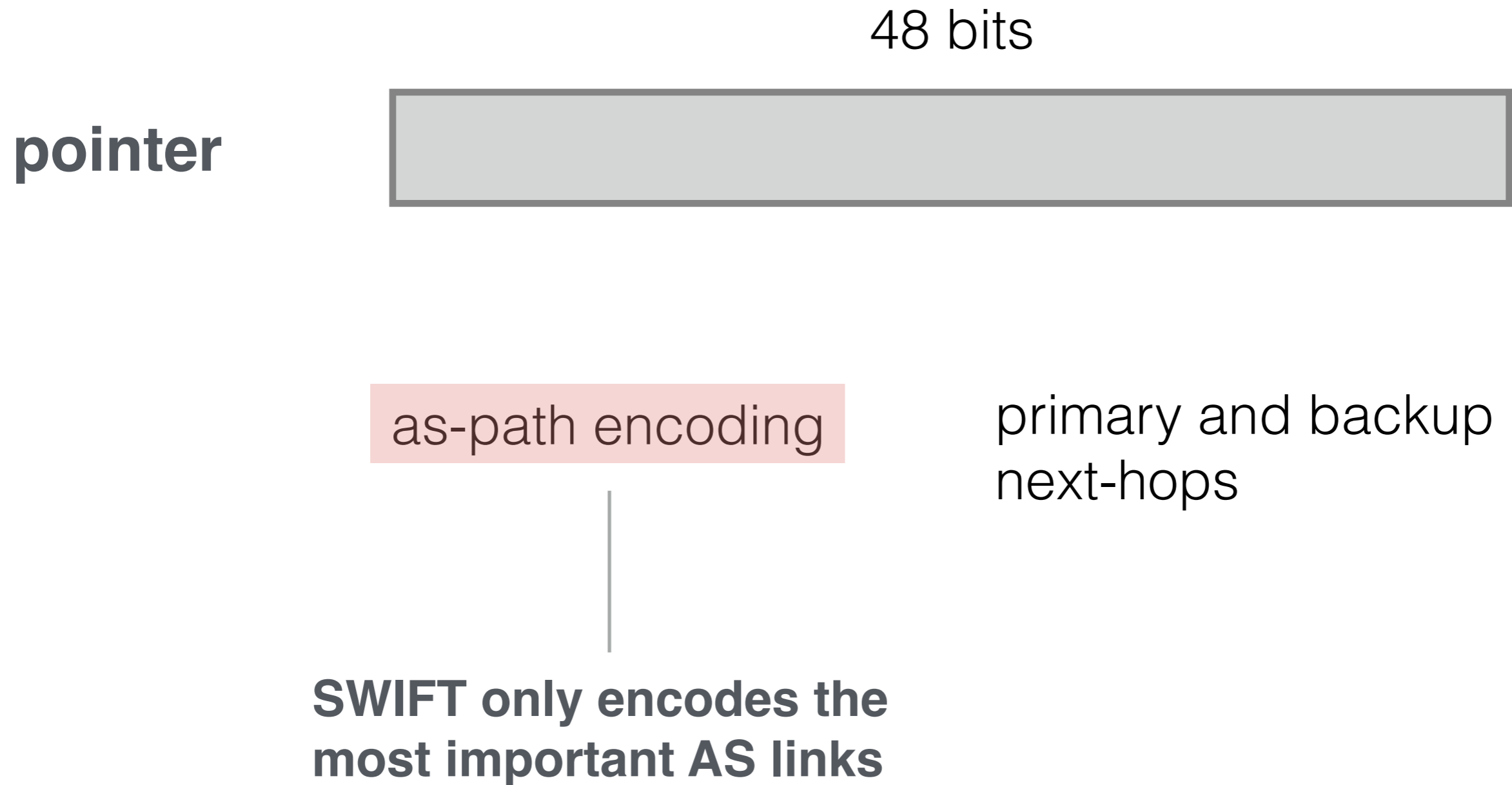
Pointers are data-plane tags with limited size



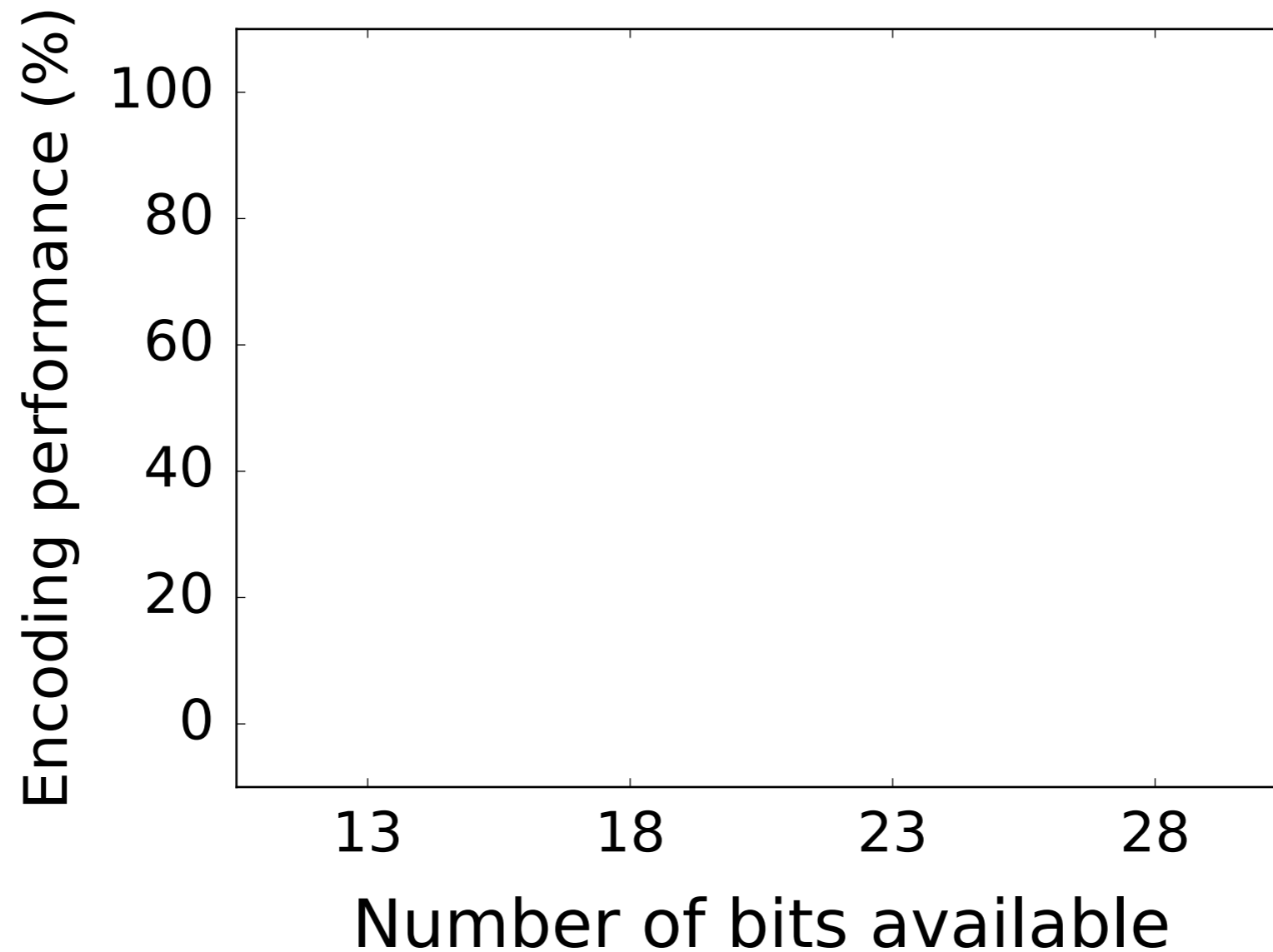
Pointers are data-plane tags with limited size



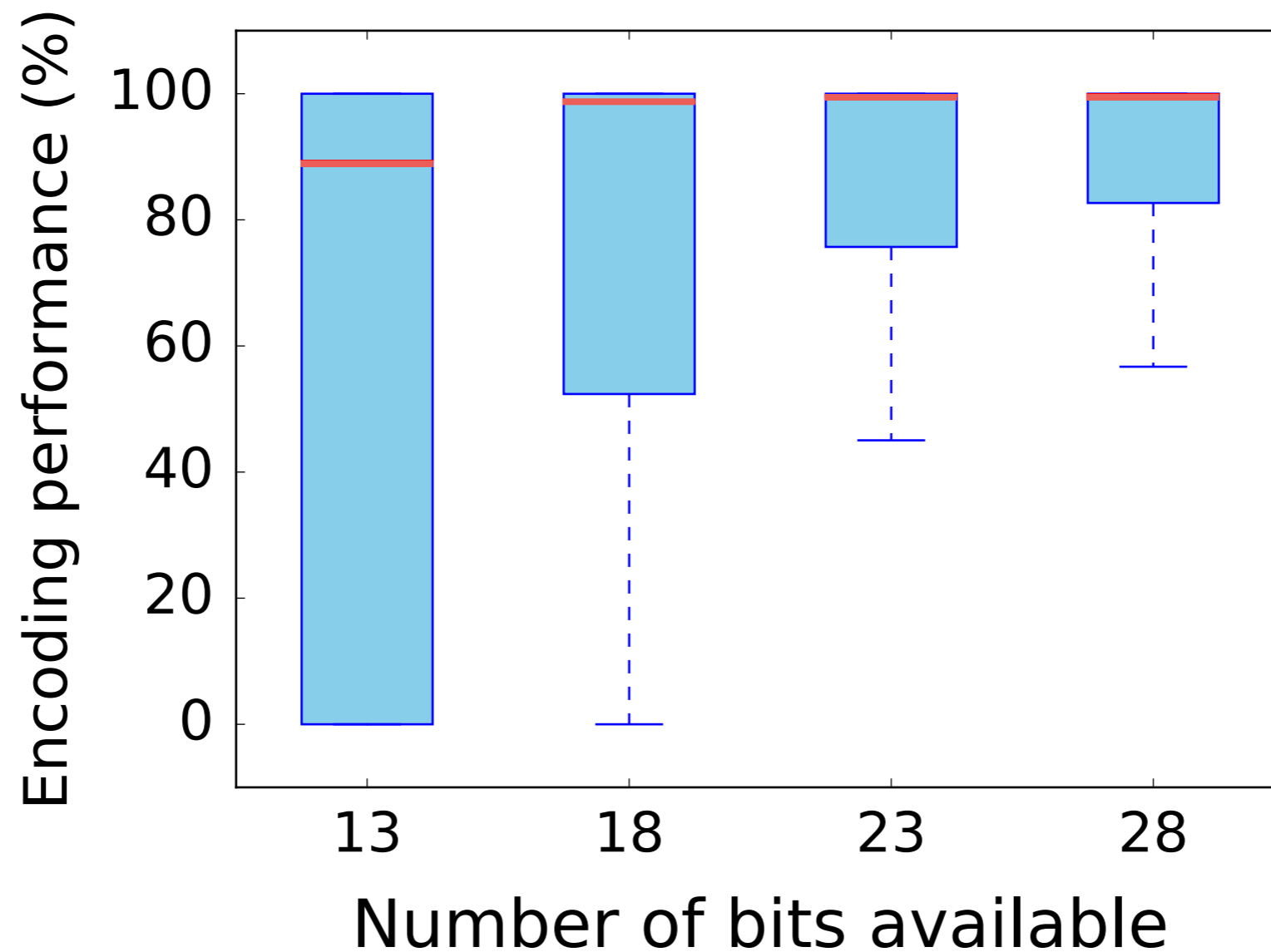
Pointers are data-plane tags with limited size



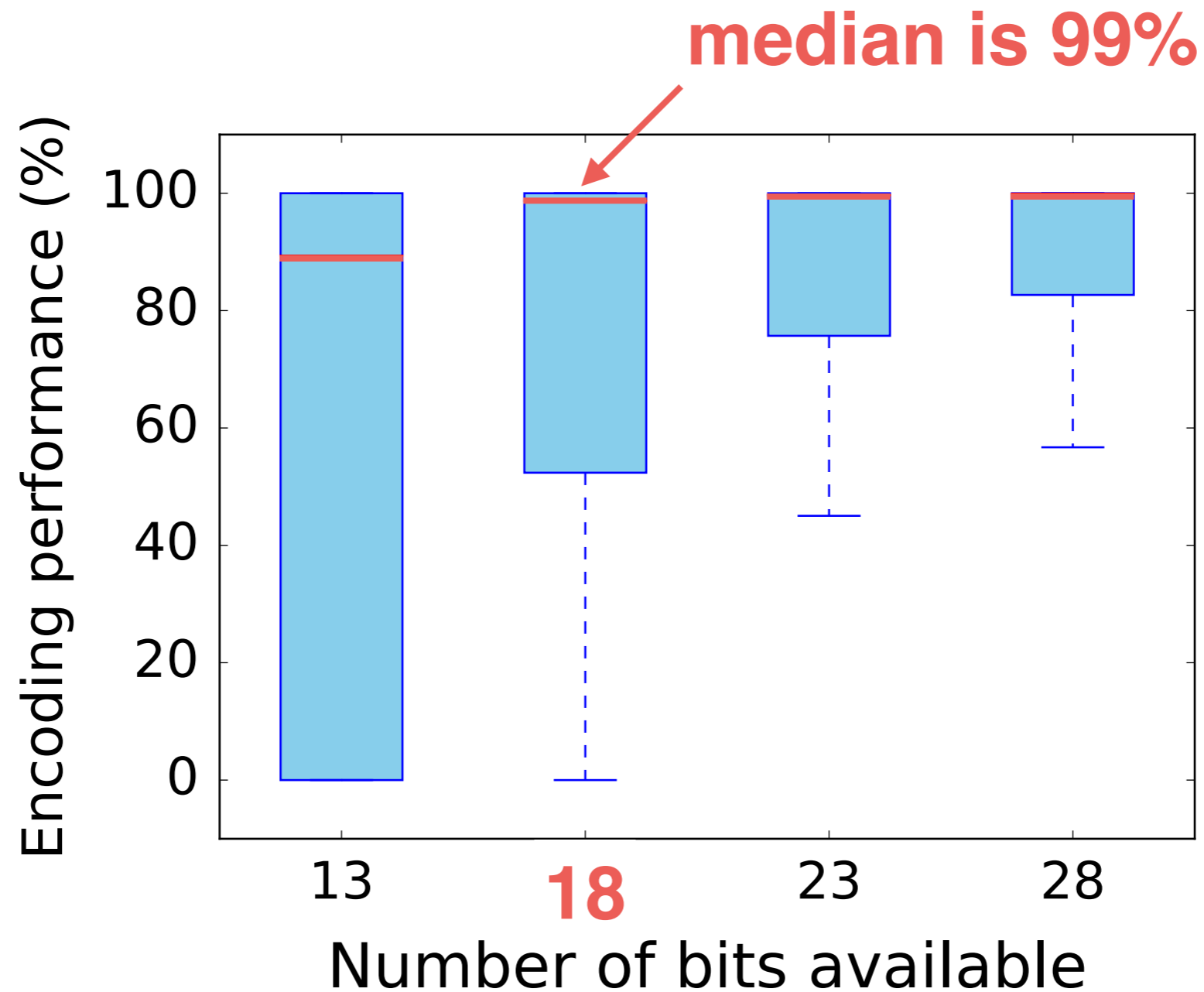
SWIFT as path encoding algorithm enables to reroute most of the affected prefixes, with few bits only



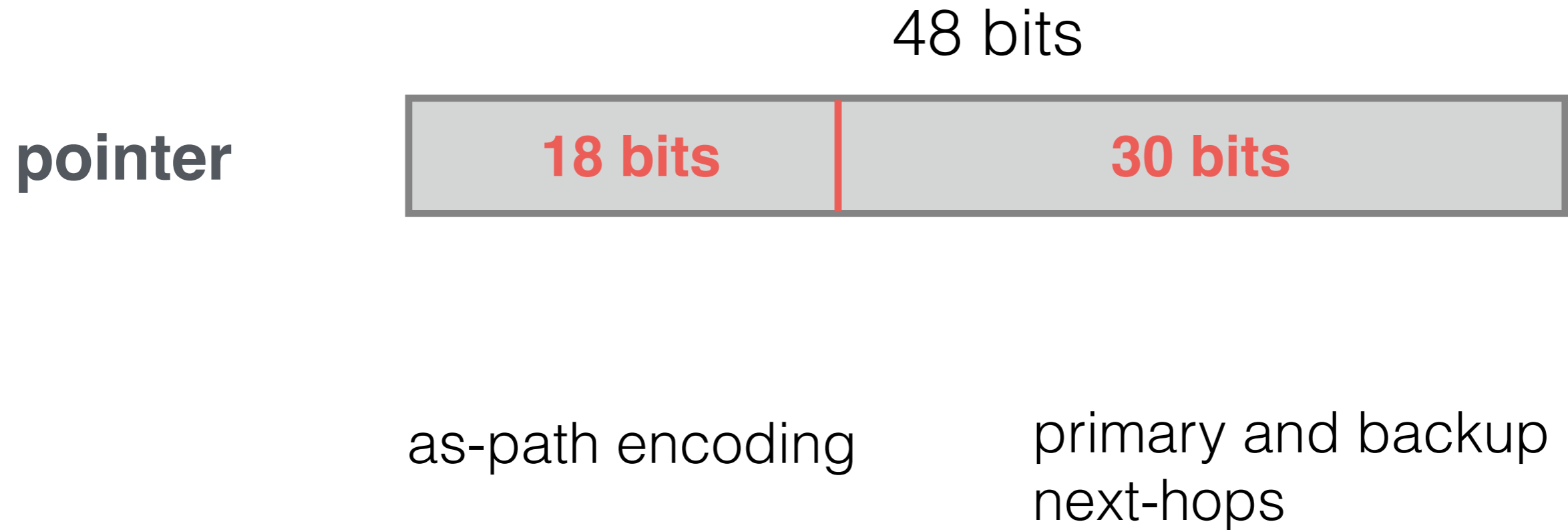
SWIFT as path encoding algorithm enables to reroute most of the affected prefixes, with few bits only



SWIFT as path encoding algorithm enables to reroute most of the affected prefixes, with few bits only



Pointers are data-plane tags with limited size



SWIFT requires less than 100 forwarding updates to reroute (with 16 backups)

SWIFT requires less than **100 forwarding updates**
to reroute (with 16 backups)

40 ms

SWIFT: Predictive Fast Reroute

1. **SWIFT** reduces the learning time of a withdrawal from **13s to 2s** using a control-plane prediction algorithm
2. **SWIFT** fast reroutes the affected traffic towards unaffected paths in **40ms** using a hierarchical forwarding table
3. **SWIFT** is **deployable** on existing routers

SWIFT: Predictive Fast Reroute

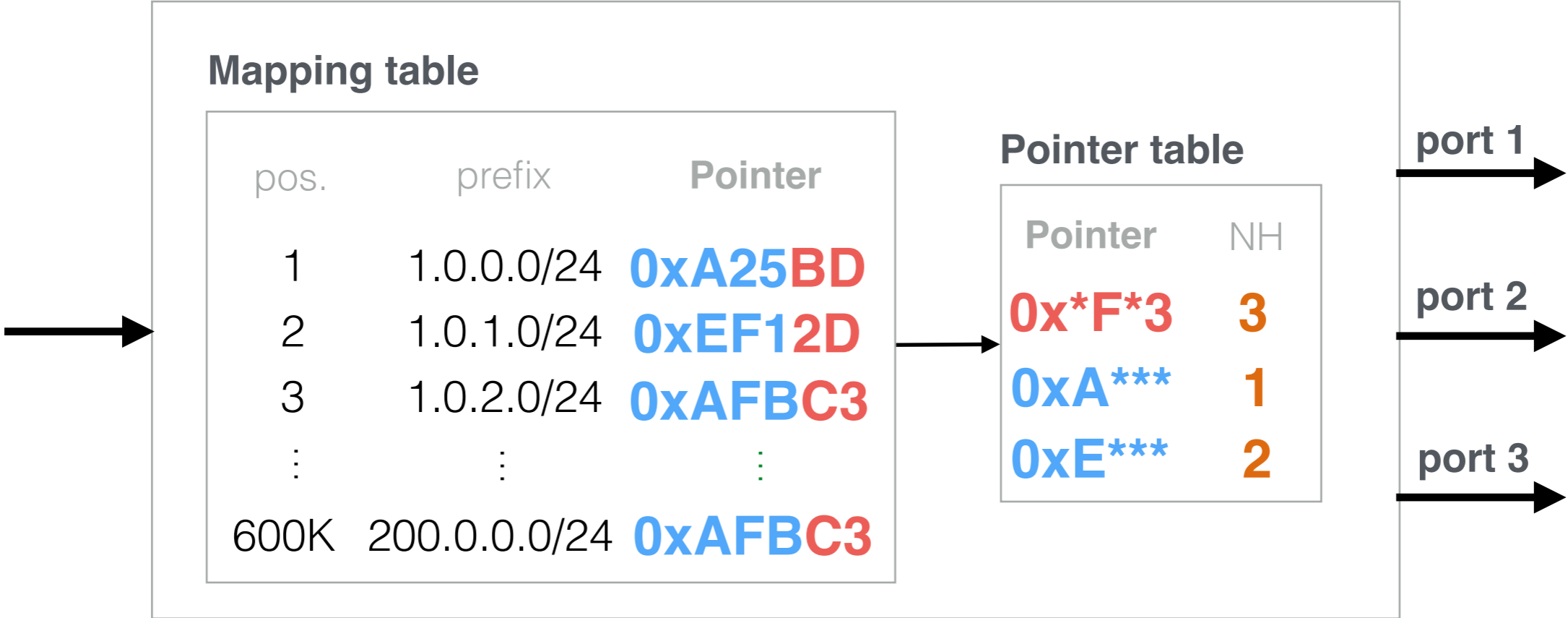
1. ***SWIFT*** reduces the learning time of a withdrawal from **13s to 2s** using a control-plane prediction algorithm
2. ***SWIFT*** fast reroutes the affected traffic towards unaffected paths in **40ms** using a hierarchical forwarding table
3. ***SWIFT*** is **deployable** on existing routers

SWIFT is deployable on existing routers equipped with a hierarchical forwarding table

SWIFT is deployable on existing routers
equipped with a hierarchical forwarding table

|
what if we don't have one though?

Hierarchical Forwarding



IP Router

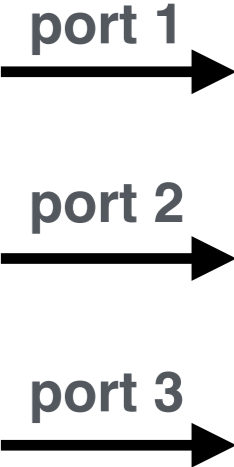
Mapping table

pos.	prefix	Pointer
1	1.0.0.0/24	0xA25BD
2	1.0.1.0/24	0xEF12D
3	1.0.2.0/24	0xAFBC3
⋮	⋮	⋮
600K	200.0.0.0/24	0xAFBC3

SDN switch

Pointer table

Pointer	NH
0x*F*3	3
0xA***	1
0xE***	2



IP Router

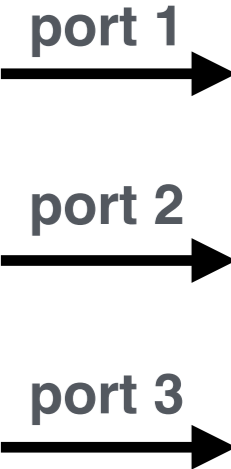
Mapping table

pos.	prefix	Dest. MAC
1	1.0.0.0/24	0xA25BD
2	1.0.1.0/24	0xEF12D
3	1.0.2.0/24	0xAFBC3
⋮	⋮	⋮
600K	200.0.0.0/24	0xAFBC3

SDN switch

Pointer table

Dest. MAC	NH
0x*F*3	3
0xA***	1
0xE***	2

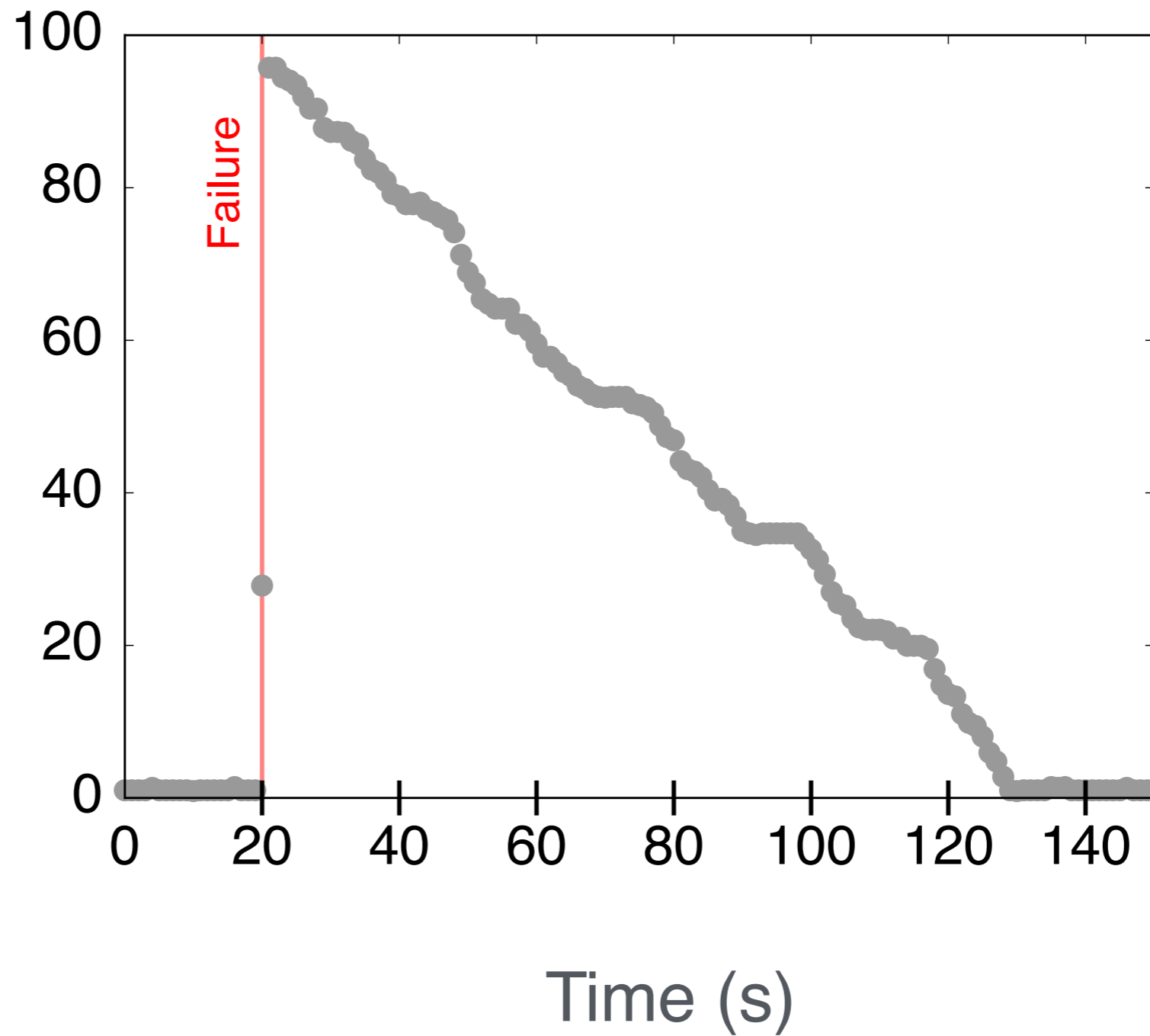


**Cisco
Router**

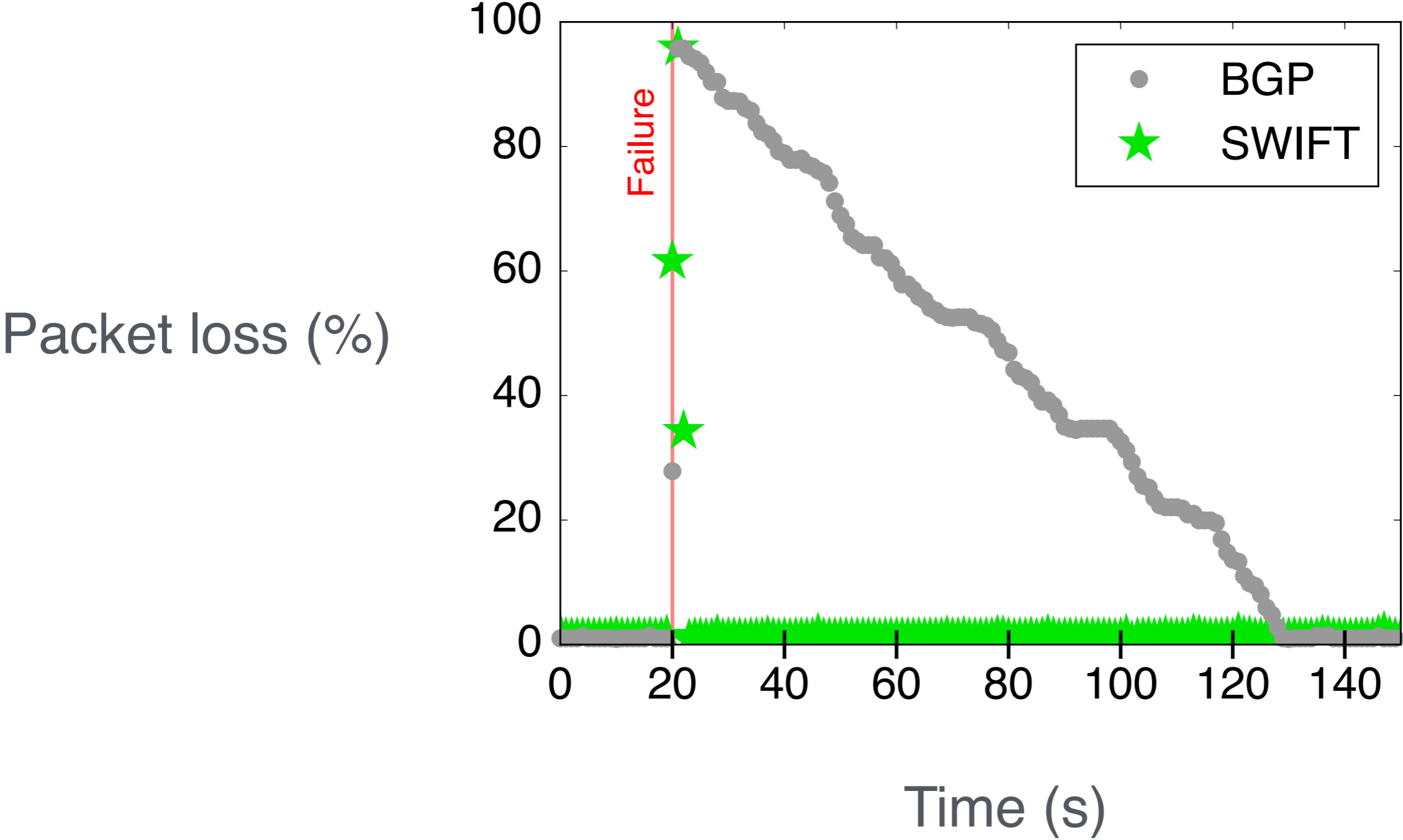
***SWIFT*
controller**

**SDN
switch**

Packet loss (%)



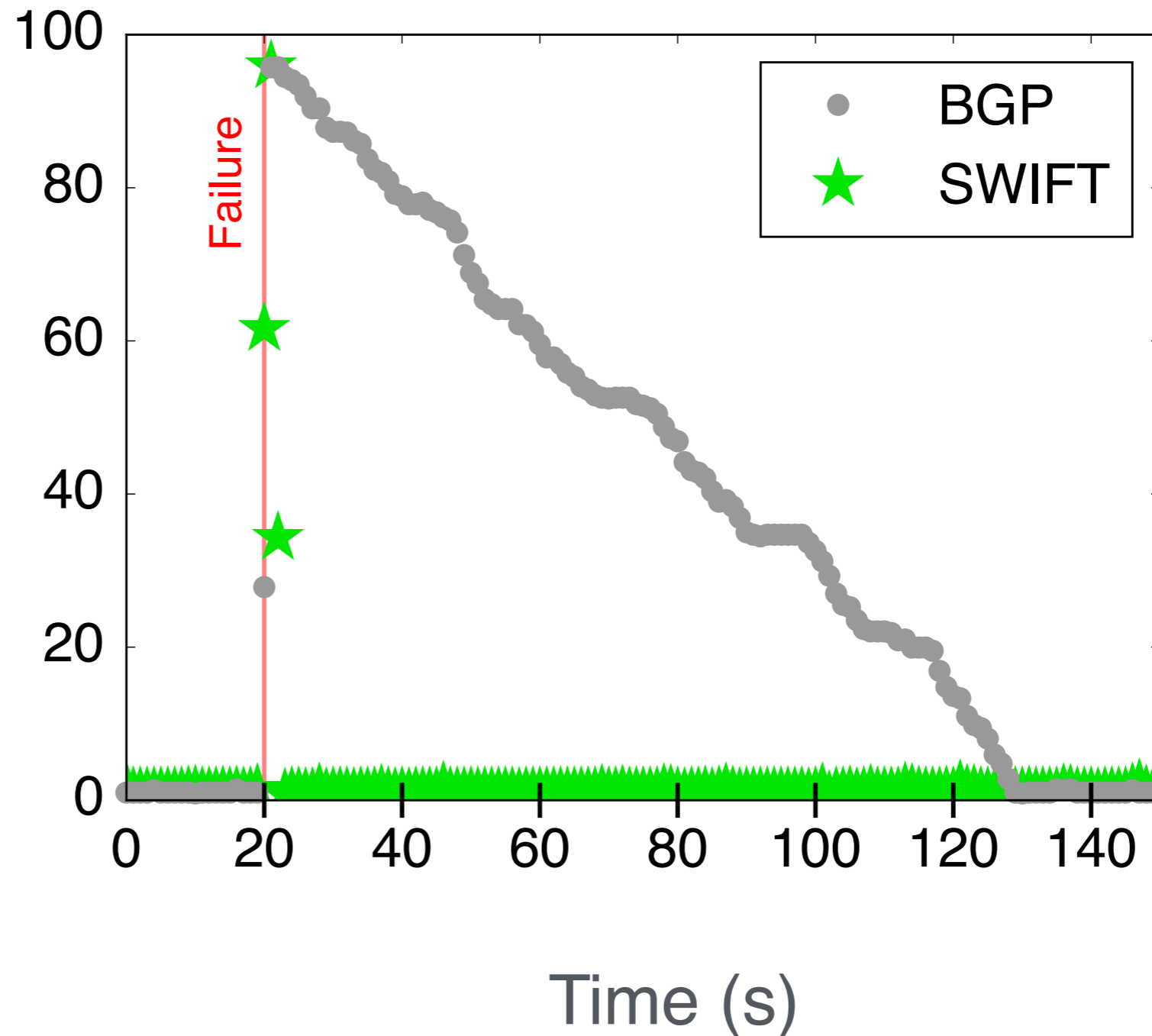
SWIFT reduces the convergence time from ~110s to 2s



SWIFT reduces the convergence time from **~110s to 2s**

**98% speed up
improvement**

Packet loss (%)



SWIFT: Predictive Fast Reroute

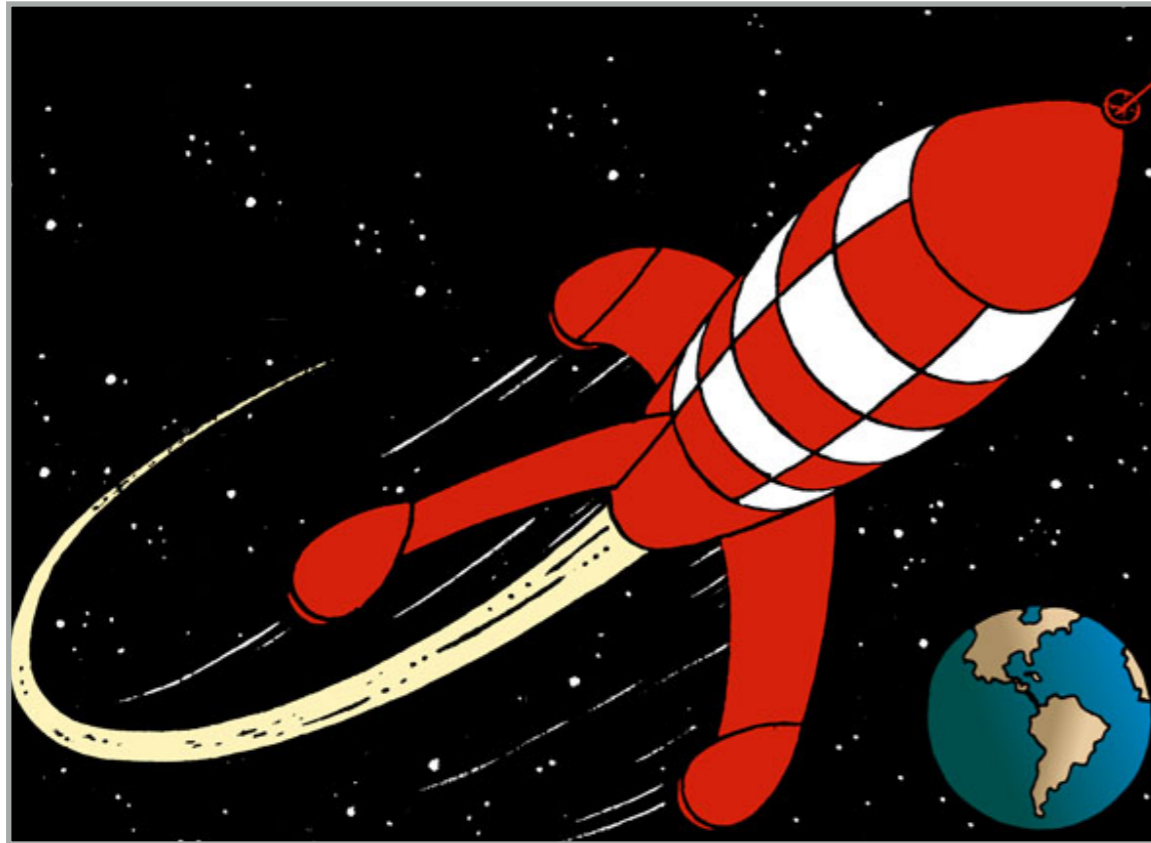
Speeds up the learning phase by predicting the control-plane
85% accuracy early in the burst

Quickly fast reroutes the affected traffic to unaffected paths
for 99% of the prefixes

Is deployable on existing routers
98% speed up improvement

swift.ethz.ch (source code, VM + demo!)

SWIFT: Predictive Fast Reroute



Thomas Holterbach

ETH Zürich / CAIDA

SIGCOMM

24th August 2017

swift.ethz.ch

Joint work with

Stefano Vissicchio

Alberto Dainotti

Laurent Vanbever

UCLondon

CAIDA, UC San Diego

ETH Zürich