## **RouteScout** Performance-Driven Internet Path Selection



# default Internet routing protocol

default Internet routing protocol BGP statically routes traffic to despite Internet's path divers

BGP statically routes traffic to each destination prefix via a single path despite Internet's path diversity, leading to sub-optimal performance.

#### Can we dynamically route traffic according to performance?

#### Can we dynamically route traffic according to performance?

## Can we dynamically route traffic according to performance? ...why now?

## Can we dynamically route traffic according to performance? ...why now?

sharper focus on reliably high network performance

### Can we dynamically route traffic according to performance? ...why now?

- sharper focus on reliably high network performance
- increasing diversity in available paths

### Can we dynamically route traffic according to performance? ...why now?

- sharper focus on reliably high network performance
- increasing diversity in available paths

programmable devices allows scalable monitoring & flexible forwarding

#### RouteScout: Performance-Driven Internet Path Selection





#### Maria Apostolaki

#### Ankit Singla





#### Laurent Vanbever

#### RouteScout: Performance-Driven Internet Path Selection

### RouteScout Operations RouteScout Architecture Data-Plane Design Evaluation

#### RouteScout: Performance-Driven Internet Path Selection

#### RouteScout Operations

#### RouteScout Architecture

#### Data-Plane Design

Evaluation

#### RouteScout selects the next hop per prefix



#### RouteScout selects the next hop per prefix



#### Routes via AS A and via AS B to AS C and AS D are equally preferred



#### BGP routes to both destinations via AS A



#### Unbeknownst to AS X Routes via AS B have much lower delay



#### What does RouteScout need to explore and exploit this opportunity?

### RouteScout runs at the network edge on a programmable switch and controls the paths of outgoing traffic



### RouteScout runs at the network edge on a programmable switch and controls the paths of outgoing traffic



#### RouteScout directs traffic to alternative BGP-compliant next-hops





dstD	loss	delay
Α	0	5
В	0	4

dstC	loss	delay
Α	0	13
В	0	7



dstD	loss	delay
Α	0	5
В	0	4

dstC	loss	delay
Α	0	13
В	0	7





### RouteScout using performance monitoring as input, solves an optimization problem to find the most suitable forwarding plan



dstD	loss	delay
Α	0	5
В	0	4

dstC	loss	delay
Α	0	13
В	0	7

### RouteScout using performance monitoring as input, solves an optimization problem to find the most suitable forwarding plan



dstD	loss	delay
Α	0	5
В	0	4

dstC	loss	delay
Α	0	13
В	0	7













#### RouteScout optimizes aggregate performance while respecting constraints





#### RouteScout optimizes aggregate performance while respecting constraints




## RouteScout optimizes aggregate performance while respecting constraints





# RouteScout Operations RouteScout Architecture Data-Plane Design Evaluation









RouteScout Operations RouteScout Architecture Data-Plane Design Loss Monitor Delay Monitor Evaluation

RouteScout Operations RouteScout Architecture Data-Plane Design Loss Monitor Delay Monitor



mirror all traffic!



mirror all traffic!



mirror all traffic!

> inflexible & nonscalable

store packets in the data plane!



store packets in the data doesn't fit in O(10)MB SRAM

mirror all traffic!





store packets in the data doesn't fit in O(10)MB SRAM

mirror all traffic!

• •

## pull Bloom Filter periodically!

stress the control-data link

> store packets in the data doesn't fit in O(10)MB SRAM

mirror all traffic!



\*Holterbach et. al, NSDI '19 Blink: Fast Connectivity Recovery Entirely in the Data Plane mirror all traffic!

inflexible & nonscalable

store packets in the data doesn't fit in O(10)MB SRAM



Blink: Fast Connectivity Recovery Entirely in the Data Plane

mirror all traffic!

inflexible & nonscalable

store packets in the data doesn't fit in O(10)MB SRAM

## The Loss Monitor... performs monitoring & aggregation in the data plane

# The Loss Monitor... performs monitoring & aggregation in the data-plane minimizes communication between the control- and data-plane

The Loss Monitor... performs monitoring & aggregation in the data plane scales with #flows (not #packets)

minimizes communication between the control and data plane

The Loss Monitor... performs monitoring & aggregation in the data plane minimizes communication between the control and data plane scales with #flows (not #packets)

## To measure loss uses a Count Min Sketch and an Aggregator

(	Coun
1	
2	
3	
4	
5	
6	
7	
8	



### Loss Aggregator

Prefix, NextHop	Expected	Lost
C,A		
С, В		

## The Loss monitor uses a Count Min Sketch to detect retransmissions

Count Min Sketch	
1	0
2	0
3	1
4	2
5	0
6	0
7	0
8	0

## Each non-empty packet, say Pkt1, triggers the insertion of the next expected one to the CMS



Count Min Sketch	
1	0
2	0
3	1
4	2
5	0
6	0
7	0
8	0

### Each non-empty packet, say Pkt1, triggers the insertion of the next expected one to the CMS Count Pkt1 IP 11.1.1.3 11.1.1.5 7 2 **TCP** 11 111 ACK 3 **SEQ #6000** LEN #500 4 Expected SEQ = 6000 + 500 = 65005 6

7

8

: Min Sketch
0
0
1
2
0
0
0
0

## Each non-empty packet, say Pkt1, triggers the insertion of the next expected one to the CMS



t Min Sketch
0+1
0
1 <b>+ 1</b>
2 <b>+1</b>
0
0
0
0

## Step 1: Pkt2 arrives, it verifies that it is expected



Count Min Sketch	
1	1
2	0
3	2
4	3
5	0
6	0
7	0
8	0

### Step 1: Pkt2 arrives, it verifies that it is expected



## Step 2: Pkt2 triggers an update in the Loss Aggregator



### Step 3: Pkt2 cleans the filters



t	Min Sketch
	1 <b>- 1</b>
	0
	2 <b>-1</b>
	3 <b>-1</b>
	0
	0
	0
	0

## Step 4: Pkt2 triggers the insertion the next expected packet, Pkt3



Min Sketch	
0	
0	
1	
2 <b>+1</b>	
0+1	
0	
0	
0+ <b>1</b>	

## If Pkt2 is lost, it will be retransmitted



Count Min Sketch	
1	0
2	0
3	1
4	3
5	1
6	0
7	0
8	1

## If Pkt2 is lost, it will be retransmitted





: Min Sketch
0
0
1
3
1
0
0
1

## The retransmitted packet, Pkt2', will find that it is not expected





## The retransmitted packet, Pkt2', will find that it is not expected



: Min Sketch
0
0
1
3
1
0
0
1

### Loss Aggregator

Prefix, NextHop	Expected	Lost
C,A	13	2 +1
С, В	0	0



RouteScout Operations RouteScout Architecture Data-Plane Design Loss Monitor Delay Monitor Evaluation
























The Delay Monitor...

does not require bidirectional traffic





- The Delay Monitor...
  - does not require bidirectional traffic
  - performs monitoring in the data plane
  - performs aggregation in the data plane



### To measure delay uses an Invertible Bloom Filter and an Aggregator

	Accumulator	CNT
1	0	0
2	0	0
3	0	0
4	0	0
5	0	0
6	0	0
7	0	0
8	0	0

Delay Aggregator		
Prefix, NextHop	Sum Delays	# Dela



### Upon arrival of a SYN packet its timestamp (T1) is XORed to multiple indexes



mulator	CNT
⊕T <u>1</u>	0
0	0
0	0
0	0
⊕T <u>1</u>	0
⊕T <u>1</u>	0
0	0
0	0



### Upon arrival of a SYN packet the corresponding counters increase by one



mulator	CNT
T <u>1</u>	0+1
0	0
0	0
0	0
T1	0+1
T1	0+1
0	0
0	0

### Similar updates happen for each SYN packets



mulator	CNT
. ⊕ <b>T</b> 2	2
0	0
0	0
0	0
T1	1
. ⊕ <b>T</b> 2	2
0	0
⊕T2	1



mulator	CNT
. ⊕ <b>T</b> 2	2
0	0
0	0
0	0
T1	1
. ⊕ <b>T</b> 2	2
0	0
Т2	1



mulator	CNT
. ⊕ <b>T</b> 2	2
0	0
0	0
0	0
T1	1
⊕ <b>T</b> 2	2
0	0
Т2	1



mulator	CNT
. ⊕ <b>T</b> 2	2
0	0
0	0
0	0
Τ1	1
⊕ <b>T</b> 2	2
0	0
Т2	1

delay=Tack-Tsyn= T2-T3



mulator	CNT
. ⊕ <b>T</b> 2	2
0	0
0	0
0	0
T1	1
⊕ <b>T</b> 2	2
0	0
Т2	1

Delay Aggregator		
Prefix, NextHop	Sum Delays	# Del
Α	+=(T3-T2)	+=
В		

delay=Tack-Tsyn=T3-T2



### By XORing the timestamp to all indexes and the counters by one the IBF is cleared



mulator	CNT
$\mathbf{T}_{2}=\mathbf{T}_{1}$	2-1
0	0
0	0
0	0
Τ1	1
$\mathbf{T}_{2}=\mathbf{T}_{1}$	2-1
0	0
⊕T2=0	1-1



mulator	CNT
T1	1
0	0
0	0
0	0
T1	1
T1	1
0	0
0	0

## RouteScout Operations RouteScout Architecture Data-Plane Design Evaluation

#### RouteScout Operations

### RouteScout Architecture

Data-Plane Design

Evaluation

Accuracy vs Memory Footprint Complexity vs Run Time Feasibility

#### RouteScout Operations

#### RouteScout Architecture

Data-Plane Design

Evaluation

Accuracy vs Memory Footprint Complexity vs Run Time Feasibility

(%)	
acy	
cur	
nac	
S – I	
Los	

10 15

Time (sec)

### 20 25 30

### Using 625KB (640K elements) and 2 hashes, the Loss monitor calculates the loss rate with high accuracy



### Using 625KB (640K elements) and 2 hashes, the Loss monitor calculates the loss rate with high accuracy



5	
>	
Í.	
0	
Ļ	
>	
· —	
0	
Ž	

Time (sec)

5

#### 10 15 20 25 30

### The probability of an ACK to decode its SYN's timestamp is >95% with a 1MB (640K elements) delay monitor of 2 hashes



Time (sec)

#### RouteScout Operations

#### Routescout Architecture

Data-Plane Design

Evaluation

Accuracy vs Memory footprint Complexity vs Run Time Feasibility



#### **#Prefixes**



### RouteScout is fast even when run with an increasing number of next hops and destinations





#### RouteScout Operations

### RouteScout Architecture

Data-Plane Design

Evaluation

Accuracy vs Memory Footprint Complexity vs Run Time Feasibility

#### RouteScout Operations

### RouteScout Architecture

Data-Plane Design

Evaluation

Accuracy vs Memory Footprint Complexity vs Run Time Feasibility

RouteScout is a modern answer to the old problem of performance-aware Internet routing.

RouteScout is a closed-loop control system with hardware and software components that leverages programmable data planes.

RouteScout's data plane fits in today's devices while its control plane runs in sub-second operating times.



