

dShark: A General, Easy to Program and Scalable Framework for Analyzing In-network Packet Traces

Da Yu[†] , Yibo Zhu[§] , Behnaz Arzani[§] , Rodrigo Fonseca[†] ,
Tianrong Zhang[§] , Karl Deng[§] , Lihua Yuan[§]

[†]Brown University



[§]Microsoft



Network reliability is critical

Cloudflare: A bad config (router rule) caused all of their edge routers to crash, taking down all of Cloudflare.

LILY HAY NEWMAN SECURITY 11.06.17 05:00 PM
HOW A TINY ERROR SHUT OFF THE INTERNET FOR PARTS OF THE US

Steam Outage: How to Monitor Data Connections

Posted by NI

Stack Overflow: A bad firewall config blocked stackexchange/stackoverflow.

Etsy: Sending multicast traffic without properly configuring switches caused an Etsy global outage.

How a typo took down the internet

Hello, operator

By Casey Newton | @CaseyNewton | Mar 2, 2017, 1:24pm EST

f t SHARE

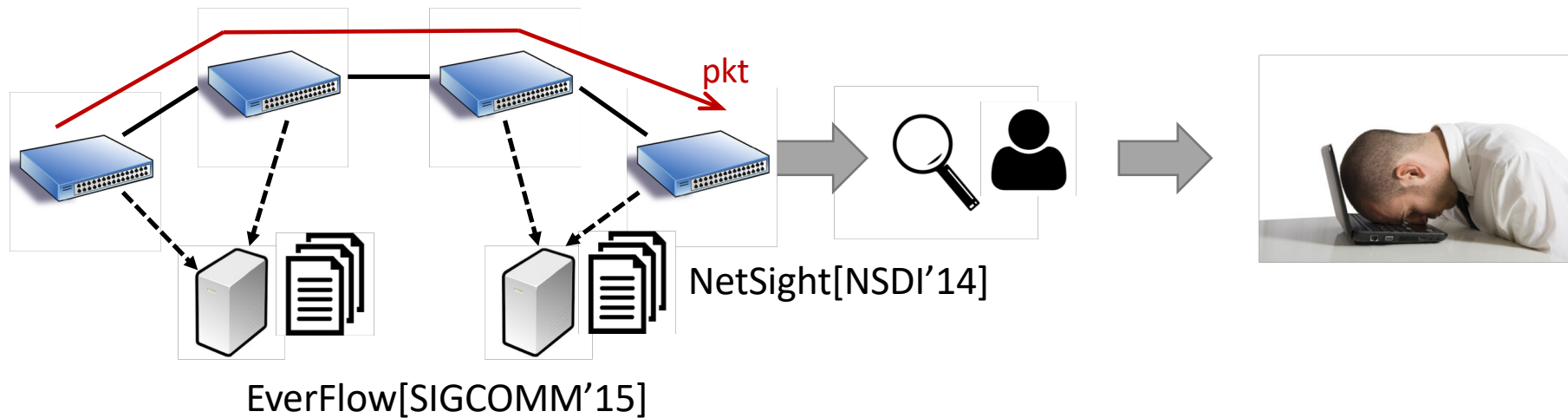
Google servers service outage due to odd traffic routing

Looks like someone pressed the wrong button on the routing machine

Existing tools are the first attempts

	End-host based	Topology or hardware specific	Target on specific problems
Trumpet[SIGCOMM'16]	X		X
Sonata[SIGCOMM'18]		X	
PathDump[OSDI'16]	X	X	
007[NSDI'18]	X	X	X
SwitchPointer[NSDI'18]	X	X	
Pingmesh[SIGCOMM'15]	X	X	X
INT		X	

In-network packet capture is the last resort



Analyzing the in-network packet traces is **challenging!**

In-network analysis: challenges

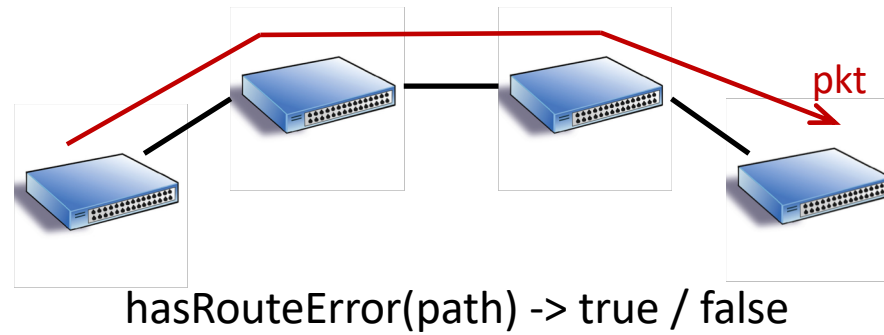
Volume

- 3.33 Mpps line-speed (10 Gbps, 1500 Bytes)

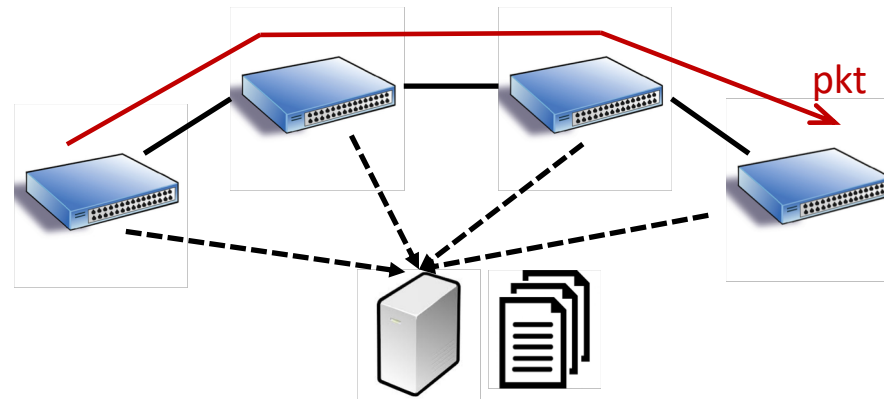
Analysis logic varies

- Logic is different case by case

Example: Route error checker

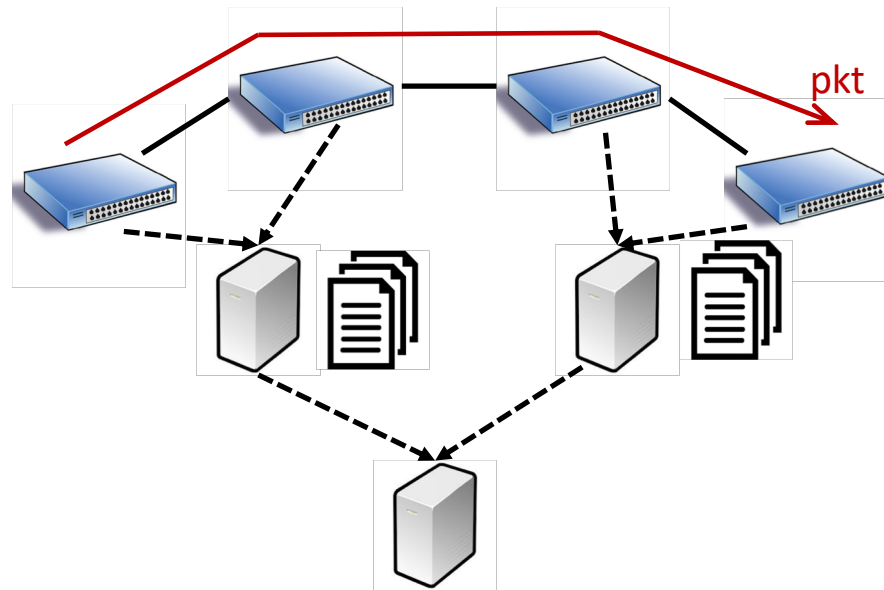


Example: Route error checker



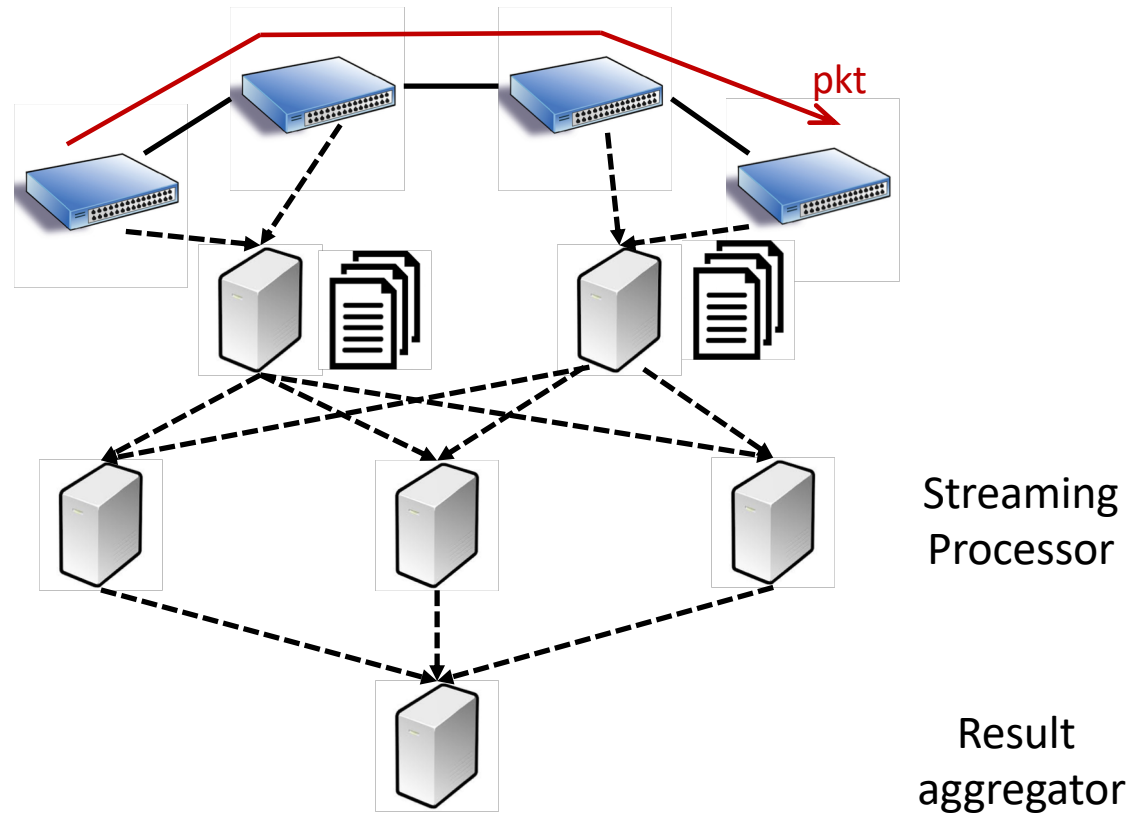
hasRouteError(path) -> true / false

Example: Route error checker



hasRouteError(path) -> true / false

Example: Route error checker



hasRouteError(path) -> true / false

In-network analysis: challenges

Volume

- 3.33 Mpps line-speed (10 Gbps, 1500 Bytes)

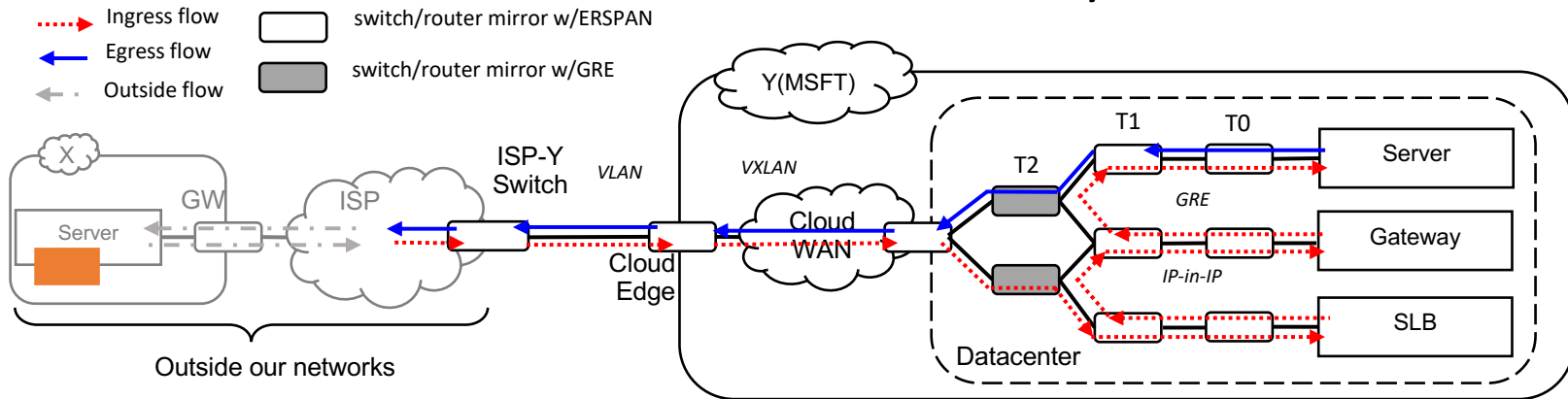
Analysis logic varies

- Logic is different case by case

Difficult to get robust analysis

- Header transformation
 - Headers are modified by the middleboxes

Packet headers are modified by middleboxes



Header format

Headers added after mirroring			Mirrored headers							
ETHERNET	IPV4	ERSPAN	ETHERNET				IPV4	TCP		
ETHERNET	IPV4	ERSPAN	ETHERNET			802.1Q	IPV4	TCP		
ETHERNET	IPV4	ERSPAN	ETHERNET	IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP	
ETHERNET	IPV4	GRE		IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP	
ETHERNET	IPV4	ERSPAN	ETHERNET		IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP
ETHERNET	IPV4	ERSPAN	ETHERNET	IPV4	IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP
ETHERNET	IPV4	GRE		IPV4	IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP
ETHERNET	IPV4	ERSPAN	ETHERNET	IPV4	IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP
ETHERNET	IPV4	ERSPAN	ETHERNET			GRE	ETHERNET	IPV4	TCP	
ETHERNET	IPV4	GRE		IPV4		GRE	ETHERNET	IPV4	TCP	
ETHERNET	IPV4	ERSPAN	ETHERNET			GRE	ETHERNET	IPV4	TCP	

Same protocol headers bring ambiguity

Header format										
Headers added after mirroring			Mirrored headers							
ETHERNET	IPV4	ERSPAN	ETHERNET						IPV4	TCP
ETHERNET	IPV4	ERSPAN	ETHERNET					802.1Q	IPV4	TCP
ETHERNET	IPV4	ERSPAN	ETHERNET		IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP
ETHERNET	IPV4	GRE			IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP
ETHERNET	IPV4	ERSPAN	ETHERNET		IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP
ETHERNET	IPV4	ERSPAN	ETHERNET	IPV4	IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP
ETHERNET	IPV4	GRE		IPV4	IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP
ETHERNET	IPV4	ERSPAN	ETHERNET	IPV4	IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP
ETHERNET	IPV4	ERSPAN	ETHERNET		IPV4		GRE	ETHERNET	IPV4	TCP
ETHERNET	IPV4	GRE			IPV4		GRE	ETHERNET	IPV4	TCP
ETHERNET	IPV4	ERSPAN	ETHERNET		IPV4		GRE	ETHERNET	IPV4	TCP

dShark: three goals

Scalable

Broadly applicable

Robust in the wild

- Header transformation

dShark: three goals

Scalable

- Components work independently and in parallel.

Broadly applicable

Robust in the wild

- Header transformation

dShark: three goals

Scalable

- Components work independently and in parallel.

Broadly applicable

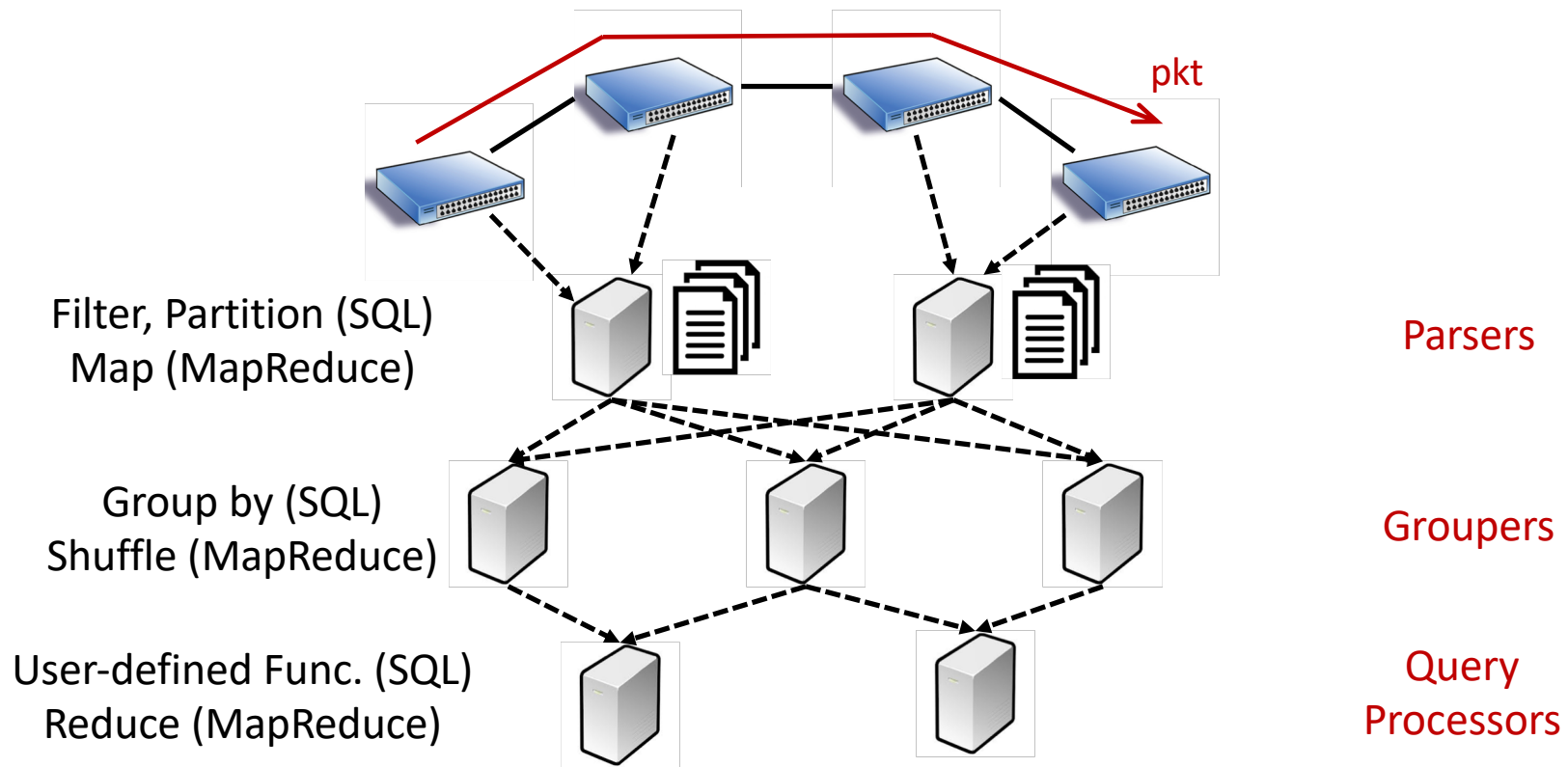
Robust in the wild

- Header transformation

How operators manually process traces

Observation #1:

- Four diagnosis steps: parse, filter, aggregate and analyze



Need to tightly integrated with the collecting infrastructure!

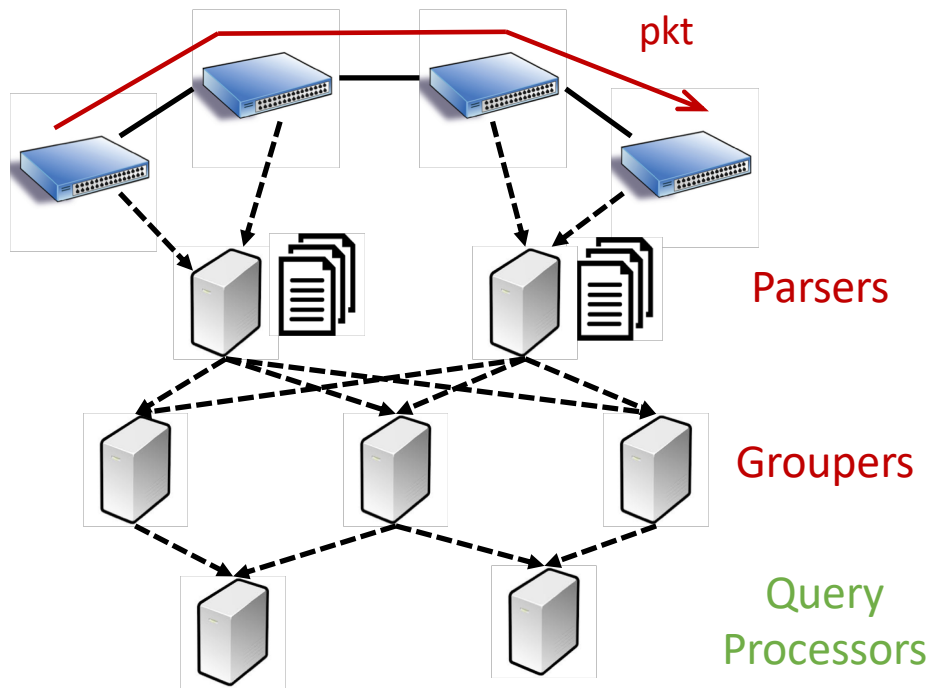
How operators manually process traces

Observation #2:

- Diagnosis logic always run on top of 4 aggregation types

	One-Hop	Multi-Hop
One-Packet	check appearance of a packet	show full path of each packet in the network
Multi-Packet	diagnose middlebox behaviors	complicated cases that requires end information

dShark's programming model



Decl. packet spec.:

- Familiar to the operators

Imp. query function:

- Flexible for analysis logic

Pkt spec.:

All instances of the same packet

Query function:
hasRouteError(path)



Declarative spec. in parsers and groupers

A packet summary is a byte array that only contains fields that the operators are interested in.

```
{  
  Summary: {  
    Key: [ipId, seqNum],  
    Additional: [ttl]  
  },  
}
```


```
  Name: {  
    ipId: ipv4.id,  
    seqNum: tcp.seq,  
    ttl: ipv4.ttl  
  }  
}
```

Definition of a packet summary


How to extract the values in the header

Same protocol headers bring ambiguity

Header format														
Headers added after mirroring			Mirrored headers											
ETHERNET	IPV4	ERSPAN	ETHERNET					IPV4	TCP					
ETHERNET	IPV4	ERSPAN	ETHERNET					802.1Q	IPV4	TCP				
ETHERNET	IPV4	ERSPAN	ETHERNET	IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP					
ETHERNET	IPV4	GRE						IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP	
ETHERNET	IPV4	ERSPAN	ETHERNET	IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP					
ETHERNET	IPV4	ERSPAN	ETHERNET	IPV4	IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP				
ETHERNET	IPV4	GRE						IPV4	IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP
ETHERNET	IPV4	ERSPAN	ETHERNET	IPV4	IPV4	UDP	VXLAN	ETHERNET	IPV4	TCP				
ETHERNET	IPV4	ERSPAN	ETHERNET	IPV4	GRE			ETHERNET	IPV4	TCP				
ETHERNET	IPV4	GRE						IPV4	GRE	ETHERNET	IPV4	TCP		
ETHERNET	IPV4	ERSPAN	ETHERNET	IPV4	GRE			ETHERNET	IPV4	TCP				



ipv4[0]



ipv4[3]
ipv4[-1]

Declarative spec. in parsers and groupers

```
{  
  Summary: {  
    Key: [ipId, seqNum],  
    Additional: [ttl]  
  },  
  
  Name: {  
    ipId: ipv4[-1].id,  
    seqNum: tcp[-1].seq,  
    ttl: ipv4[:].ttl  
  }  
}
```

Definition of a packet summary

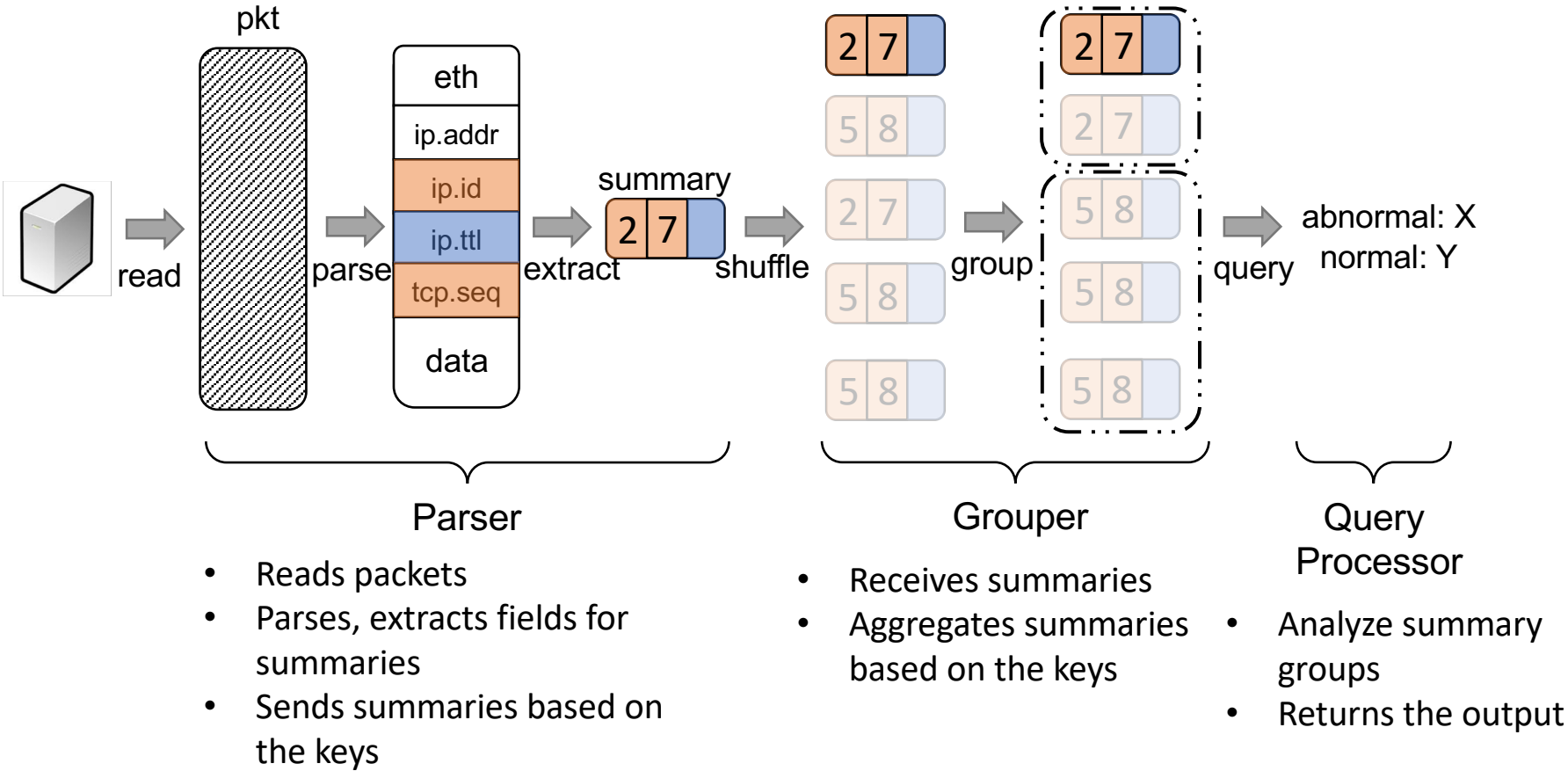
How to extract the
values in the header

Imperative diagnosis logic in the query processors

```
Pair<Object, Object> query (const vector<Summary>& group) {  
  
    // Reconstructs the path based on TTL  
    constructPath(group);  
  
    // Checks path  
    bool result = hasRouteError(group);  
  
    return make_pair(result, 1);  
}
```

- In practice, this is implemented in 49 lines of code for the query function!

dShark overview



Another example: load balancer profiler

Spec:

```
{  
  Summary: {  
    Key: [ipId, seqNum],  
    Additional: [vip, pip]  
  },  
}
```

- Innermost ipId and tcp seq # to identify a packet
- Virtual IP and the physical IP

```
Name: {  
  ipId: ipv4[-1].id,  
  seqNum: tcp[-1].seq,  
  vip: ipv4[-1].dst,  
  pip: ipv4[0].dst  
}  
}
```

Func.:

```
Pair<Pair<IP, IP>, int> query(group) {  
  // Validate data  
  ...  
  return Pair((Pair(vip, pip), 1));  
}
```

- Returns the map

- In practice, this is implemented in **18** lines of code for the query function!

Group pattern	Application	Analysis logic	In-nw ck. only	Header transf.	Query LOC
One packet on one hop	Loop-free detection [21] <i>Detect forwarding loop</i>	<i>Group:</i> same packet(ipv4[0].ipid, tcp[0].seq) on one hop <i>Query:</i> does the same packet appear multiple times on the same hop	No	No	8
	Overloop-free detection [69] <i>Detect forwarding loop involving tunnels</i>	<i>Group:</i> same packet(ipv4[0].ipid, tcp[0].seq) on tunnel endpoints <i>Query:</i> does the same packet appear multiple times on the same endpoint	Yes	Yes	8
One packet on multiple hops	Route detour checker <i>Check packet's route in failure case</i>	<i>Group:</i> same packet(ipv4[-1].ipid, tcp[-1].seq) on all switches <i>Query:</i> is valid detour in the recovered path(ipv4[:].ttl)	No	Yes*	49
	Route error <i>Detect wrong packet forwarding</i>	<i>Group:</i> same packet(ipv4[-1].ipid, tcp[-1].seq) on all switches <i>Query:</i> get last correct hop in the recovered path(ipv4[:].ttl)	No*	Yes*	49
	Netsight [21] <i>Log packet's in-network lifecycle</i>	<i>Group:</i> same packet(ipv4[-1].ipid, tcp[-1].seq) on all switches <i>Query:</i> recover path(ipv4[:].ttl)	No*	Yes*	47
	Hop counter [21] <i>Count packet's hop</i>	<i>Group:</i> same packet(ipv4[-1].ipid, tcp[-1].seq) on all switches <i>Query:</i> count record	No*	Yes*	6
Multiple packets on one hop	Traffic isolation checker [21] <i>Check whether hosts are allowed to talk</i>	<i>Group:</i> all packets at dst ToR(SWITCH=dst.ToR) <i>Query:</i> have prohibited host(ipv4[0].src)	No	No	11
	Middlebox(SLB, GW, etc) profiler <i>Check correctness/performance of middleboxes</i>	<i>Group:</i> same packet(ipv4[-1].ipid, tcp[-1].seq) pre/post middlebox <i>Query:</i> is middlebox correct(related fields)	Yes	Yes	18 [†]
	Packet drops on middleboxes <i>Check packet drops in middleboxes</i>	<i>Group:</i> same packet(ipv4[-1].ipid, tcp[-1].seq) pre/post middlebox <i>Query:</i> exist ingress and egress trace	Yes	Yes	8
	Protocol bugs checker(BGP, RDMA, etc) [69] <i>Identify wrong implementation of protocols</i>	<i>Group:</i> all BGP packets at target switch(SWITCH=tar_SW) <i>Query:</i> correctness(related fields) of BGP(FLTR: tcp[-1].src dst=179)	Yes	Yes*	23 [‡]
	Incorrect packet modification [21] <i>Check packets' header modification</i>	<i>Group:</i> same packet(ipv4[-1].ipid, tcp[-1].seq) pre/post the modifier <i>Query:</i> is modification correct (related fields)	Yes	Yes*	4 [°]
	Waypoint routing checker [21, 43] <i>Make sure packets (not) pass a waypoint</i>	<i>Group:</i> all packets at waypoint switch(SWITCH=waypoint) <i>Query:</i> contain flow(ipv4[-1].src+dst, tcp[-1].src+dst) should (not) pass	Yes	No	11
Multiple packets on multiple hops	DDoS diagnosis [43] <i>Localize DDoS attack based on statistics</i>	<i>Group:</i> all packets at victim's ToR(SWITCH=vic.ToR) <i>Query:</i> statistics of flow(ipv4[-1].src+dst, tcp[-1].src+dst)	No	Yes*	18
	Congested link diagnosis [43] <i>Find flows using congested links</i>	<i>Group:</i> all packets(ipv4[-1].ipid, tcp[-1].seq) pass congested link <i>Query:</i> list of flows(ipv4[-1].src+dst, tcp[-1].src+dst)	No*	Yes*	14
	Silent black hole localizer [43, 69] <i>Localize switches that drop all packets</i>	<i>Group:</i> packets with duplicate TCP(ipv4[-1].ipid, tcp[-1].seq) <i>Query:</i> get dropped hop in the recovered path(ipv4[:].ttl)	No*	Yes*	52
	Silent packet drop localizer [69] <i>Localize random packet drops</i>	<i>Group:</i> packets with duplicate TCP(ipv4[-1].ipid, tcp[-1].seq) <i>Query:</i> get dropped hop in the recovered path(ipv4[:].ttl)	No*	Yes*	52
	ECMP profiler [69] <i>Profile flow distribution on ECMP paths</i>	<i>Group:</i> all packets at ECMP ingress switches(SWITCH in ECMP) <i>Query:</i> statistics of flow(ipv4[-1].src+dst, tcp[-1].src+dst)	No*	No	18
	Traffic matrix [43] <i>Traffic volume between given switch pairs</i>	<i>Group:</i> all packets at given two switches(SWITCH in tar_SW) <i>Query:</i> total volume of overlapped flow(ipv4[-1].src+dst, tcp[-1].src+dst)	No*	Yes*	21

Please check the paper for details!

dShark: three goals

Scalable

- Components work independently and in parallel.

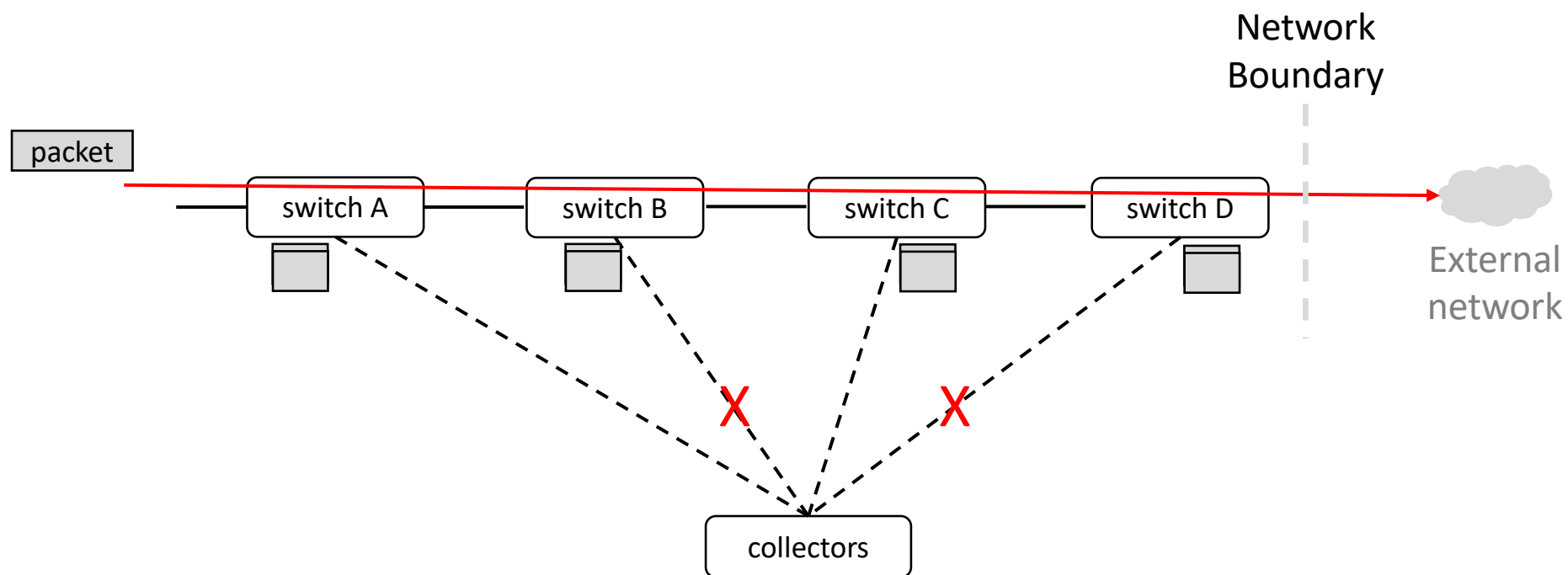
Broadly applicable

- dShark's programming model is general.
- It supports 4 types of aggregation that covers 18 typical analysis apps.

Robust in the wild

- Header transformation
 - Define packet signature in the summary.
 - Leverage the index of protocol.
- Capture noise

Tolerate capture noise



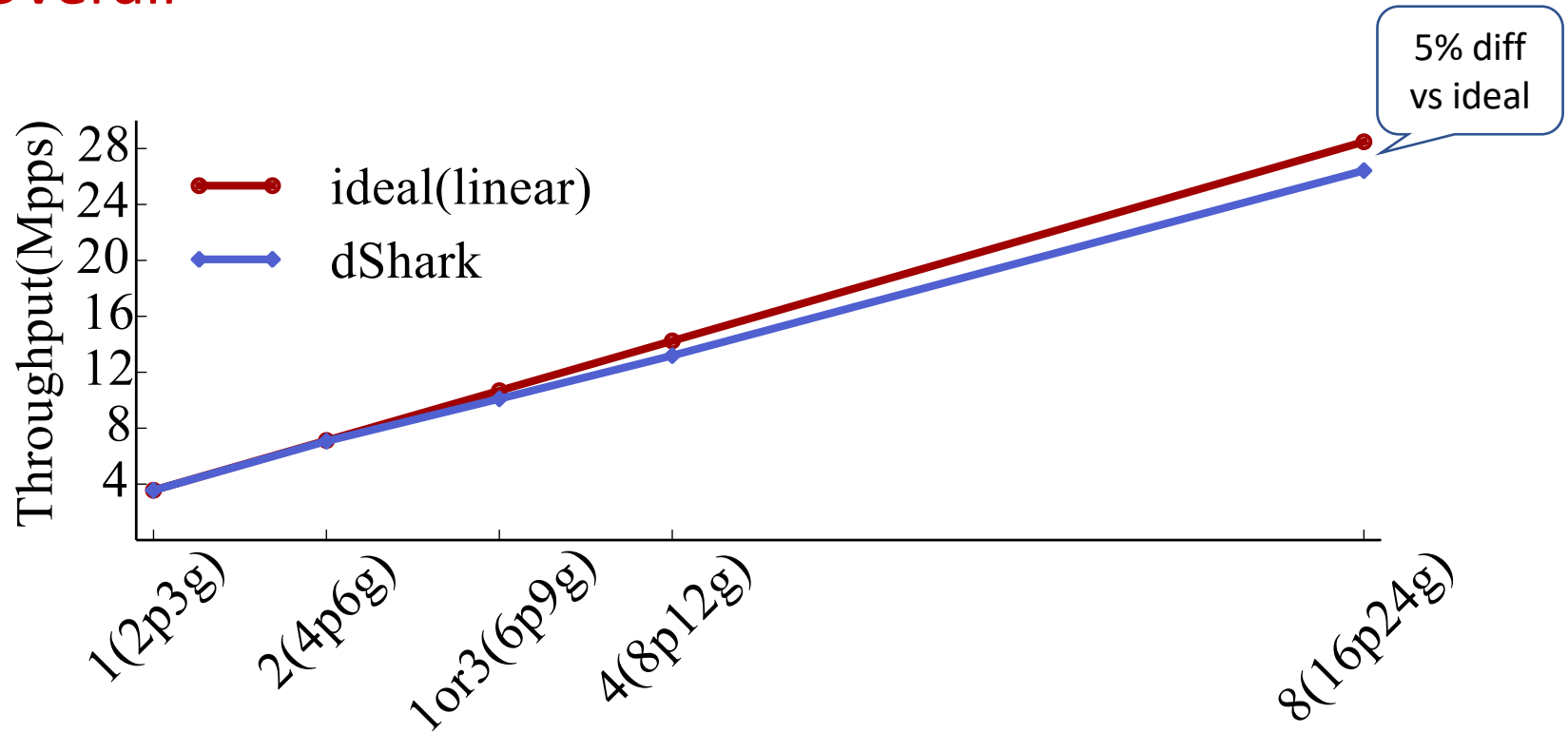
1. Recover by the next hop(s)
2. Leverage end-to-end information

Performance of dShark

- 8 VMs from a public cloud
- Each has:
 - 16-core 2.4GHz vCPU
 - 56GB memory
 - 10Gbps virtual network
- Feed with **real traces** from production

dShark scales nearly linearly

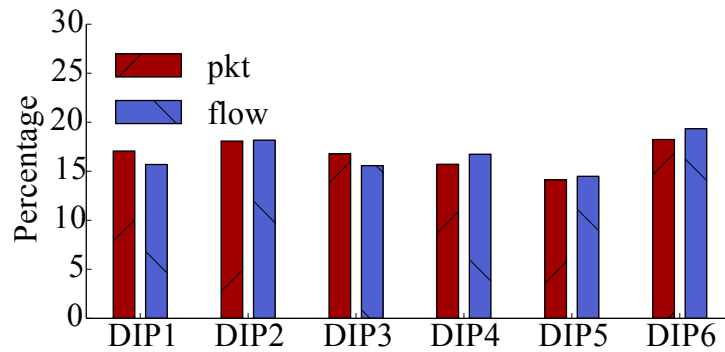
Overall



Some findings

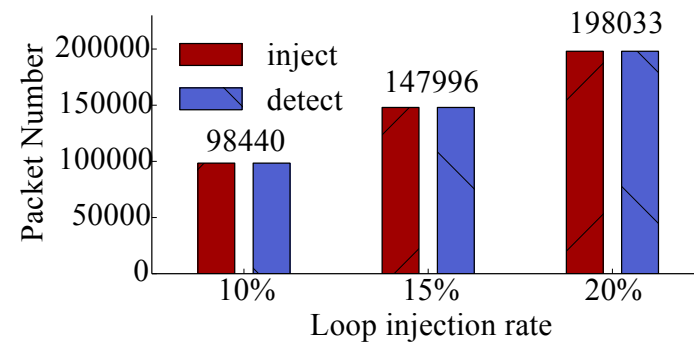
Case 1:

Profile an SLB



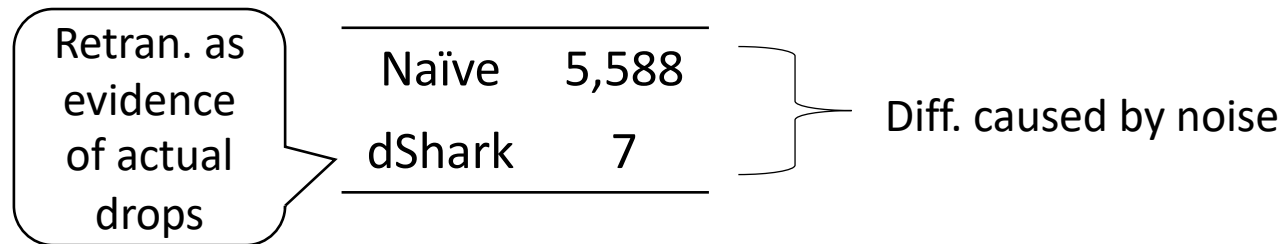
Case 2:

Detect loops



Case 3:

Detect packet drops on T1



Please check the paper for details!!

Conclusion

- dShark is a general, easy-to-program, scalable and high-performance in-network packet trace analyzer.
- Takeaways:
 - dShark's programming model is broadly applicable
 - We use this model to implement 18 different typical diagnosis apps
 - Operators focus on the logic without worrying about:
 - Header transformation, capture noise, scalability
 - dShark is fast and can scale linearly