# Programmable Per-Packet Network Telemetry: From Wire to Kafka at Scale

Zhang Liu (University of Colorado Boulder)
Yatish Kumar (LBNL, ESnet)
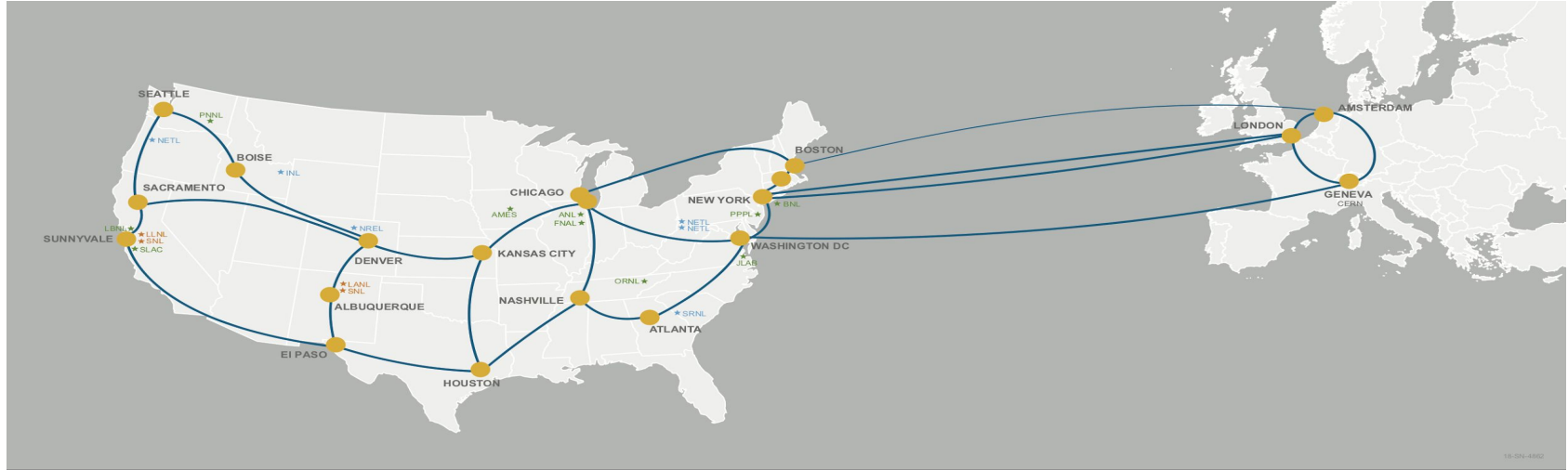Bruce Mah (LBNL, ESnet)
Chin Guok (LBNL, ESnet)
**Richard Cziva (LBNL, ESnet)**

SNTA '21, June 21, 2021, Virtual Event, Sweden

U.S. DEPARTMENT OF **ENERGY**
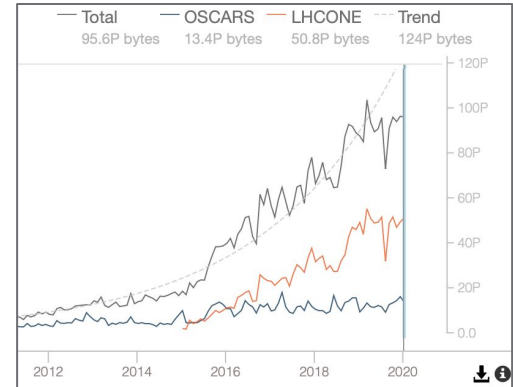Office of Science

BERKELEY LAB

# ESnet: DOE's <u>high-performance network</u> (HPN) user facility optimized for enabling big-data science



ESnet provides connectivity to
<u>all of the DOE labs</u>, experiment sites, & supercomputers

# Increasing Need for Programmability

- ESnet's traffic, user-base and the experiments continue to grow in a fast pace
- Computing and data model are also evolving, requiring:
  - fine-grained visibility in real-time
  - application-specific traffic handling
  - programmable, in-network services
- Needs not addressed by existing measurement mechanisms (sampled, aggregated, delayed)
- High Touch Services created to fulfill these needs



Live ESnet usage statistics: my.es.net
Total carried: Exabyte/year.

# Hightouch Server Hardware



2x100G
Ethernet

[4] Xilinx U280 FPGA Card

2x100G Ethernet QSFP-28

Custom Logic for Flow Tracking

High End Server

Dual Socket, Fast Storage

Hosts Hightouch Application

ESnet

# ESnet Network Packet Telemetry Data

- SNMP
  - All interfaces, 30 seconds poll interval
  - Primary use: failure detection, traffic visualization: http://my.es.net
  - Data rate: 4000 interfaces => **130 events per second**
- Netflow / IPFIX
  - All interfaces, packets sampled 1:1000
  - Primary use: capacity planning (offline)
  - Raw data rate: **~ 6500 events per second**
- High Touch Services
  - Selected interfaces and flows, 1:1 packet to telemetry
  - Primary use: high-precision telemetry
  - Raw data rate: **~ 1 to 8 million events per second** for a single interface

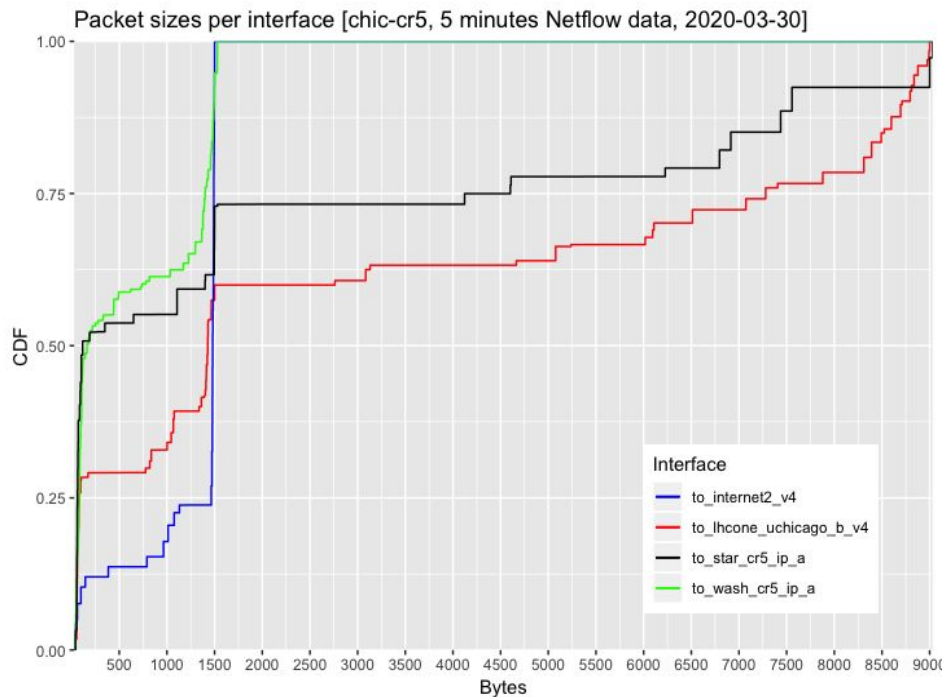| Telemetry | Raw Data Rate Per Second |
|---|---|
| SNMP | 130 |
| Netflow / IPFIX | 6500 |
| High Touch Services | 1 - 8 M |

*Telemetry Data Rates*

ESnet

# Per-Packet Data Rates

- Packet size depends on:
  - MTU
  - Application (science vs http)
  - Average for science traffic: ~1500B
- Traffic rate at ESnet at any time:
  - All traffic: O(1Tbit/s)
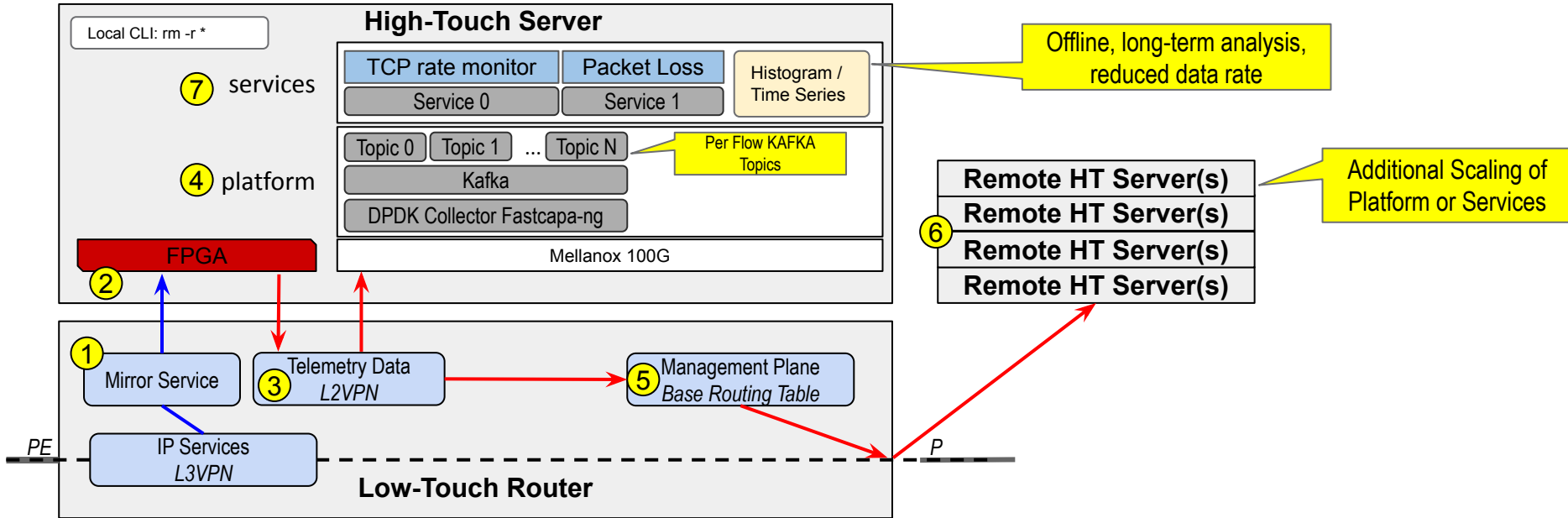  - Large customers: O(100Gbit/s)

| Packet size | Rate | Telemetry PPS | Telemetry Rate |
|---|---|---|---|
| 1500B | 10Gb/s | 812K | 1,079Mb/s |
| 1500B | 100Gb/s | 8,127K | 10,790Mb/s |
| 9000B | 10Gb/s | 138K | 183Mb/s |
| 9000B | 100Gb/s | 1,383K | 1,833Mb/s |

*Telemetry Packet Rates*



Packet sizes per interface [chic-cr5, 5 minutes Netflow data, 2020-03-30]

Interface
- to_internet2_v4
- to_lhcone_uchicago_b_v4
- to_star_cr5_ip_a
- to_wash_cr5_ip_a

*Estimated packet sizes in production*
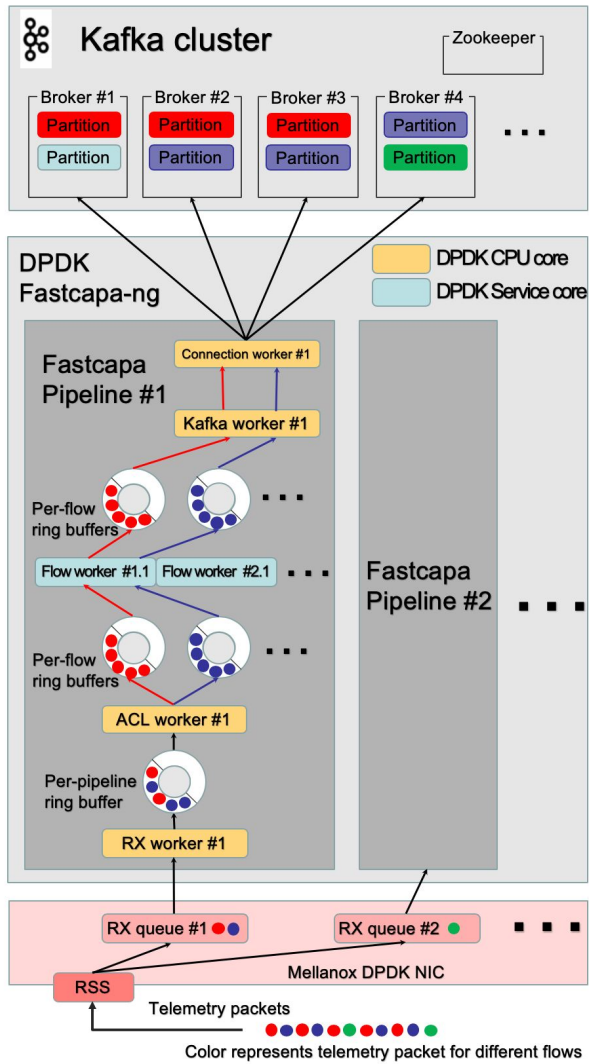
ESnet

# ESnet6 High-Touch Architecture Overview



1. Mirror Service - Allows selective flows in the dataplane to be duplicated and sent to the FPGA for processing.
2. Programmable Dataplane (DP) - Appends meta-data, timestamps and repackages packet for transmission to Platform code.
3. Telemetry Data L2VPN - Connect Dataplane and Platform, possibly on different High-Touch Servers.
4. Platform - Reads telemetry packets from the network and distributes information to High Touch Services.
5. Management Plane Base Routing Table - Provides connectivity to Remote Servers.
6. Remote Server - Hosts Platform components or Services (but not a Dataplane). Telemetry data can be directed to Remote Servers.
7. Service - Reads data from the Platform and performs real-time analysis as well as inserts selected telemetry data into database.
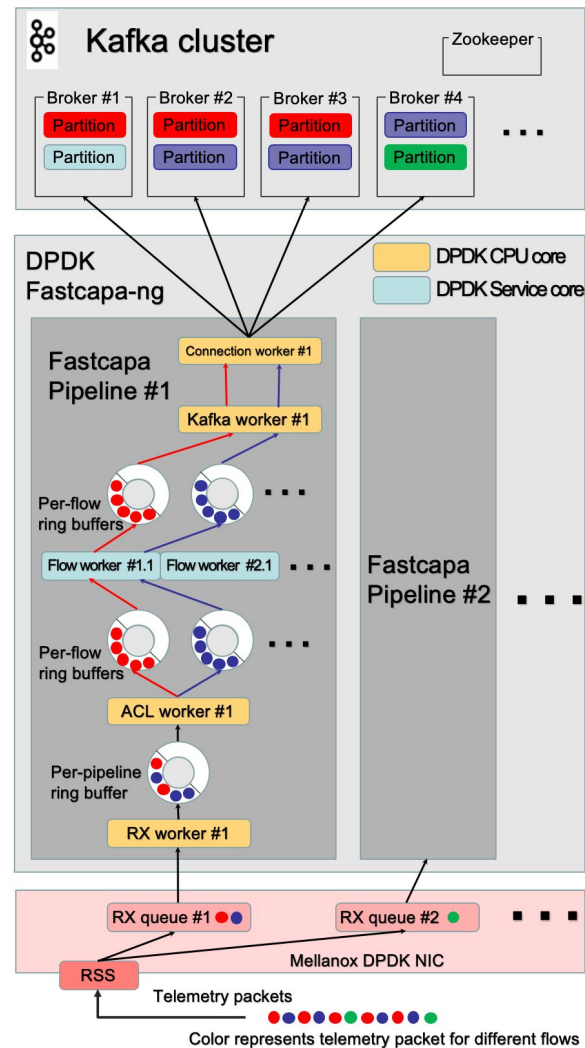
# Fastcapa-ng

- **ESnet-developed software (C / DPDK)**
  - Based on [Apache Metron Fastcapa](#)
  - Uses DPDK: fast packet processing API
  - Primary functions: telemetry processing, batching, filtering, aggregation, forwarding
- **Design goals:**
  - Packet order preservation
  - High-performance Kafka handling
  - Easy programming
- **Multi-pipeline** design for scalability, each pipeline can handle TCP flows from single 100G link.
- **Multi-stage** design for performance, each packets will be processed by 5 CPUs in series.

8

# Fastcapa-ng Internals

- Dedicated Kafka connection
  - maintain TCP connection, message compression task
- Kafka worker
  - Combine multiple telemetry packets into large kafka messages
- Flow worker (service cores)
  - process flows using different function:
    - Passthrough
    - Sampling
    - Histogram
    - (more under development)
  - Flexible N to M mapping of flow to service cores.
- ACL worker
  - classify flows and send them to dedicated rings.
- RX worker
  - pull packet into ring buffers
- RX queue
  - NIC dma packets into memory
  - RSS (Receive Side Scaling) applied

9



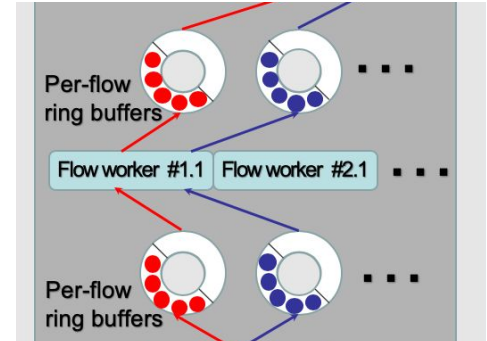Color represents telemetry packet for different flows

# Flow Worker (Service Cores)



```
108   /**
109    * This is sampling worker for service-core function.
110    * A sampling rate is defined for the flow processed by this worker.
111    */
112   static int sampling_worker(void *args)
113   {
114       service_params *params = (service_params*) args;
115       unsigned nb_in, i;
116       unsigned nb_out= 0;
117       const unsigned int flow_burst_size = params->flow_burst_size;
118       struct rte_ring *input_ring = params->input_rings[params->ring_id];
119       struct rte_ring *output_ring = params->output_rings[params->ring_id];
120
121       // dequeue packets from the ring
122       struct rte_mbuf* pkts[flow_burst_size];
123           nb_in = rte_ring_dequeue_burst(input_ring, (void*) pkts, flow_burst_size, NULL);   ←——  Read from input queue
124
125       if(likely(nb_in > 0)) {
126           params->stats.in += nb_in;
127
128           for(i = 0; i < nb_in; i ++){
129               if(params->sampling_counter == 0){
130                   rte_ring_enqueue(output_ring,pkts[i]);   ←——  Write to output queue
131                   nb_out ++;
132               }
133               else{
134                   rte_pktmbuf_free(pkts[i]);   ←——  Drop packet
135               }
136               params->sampling_counter = (params->sampling_counter + 1) % params->sampling_rate;
137           }
138           params->stats.out += nb_out;
139       }
140
141       return 0;
142   }
```

# Fastcapa-ng Runtime Configuration

```
{
    protocol = 6;
    protocol_mask = 255;
    srcIP = "192.168.25.5";
    srcIP_mask = 32;
    dstIP = "192.168.25.4";
    dstIP_mask = 32;
    srcPort = 5201;
    srcPort_mask = 5201;
    dstPort = 10001;
    dstPort_mask = 10001;
    priority = 103;

    flow_id = 3;
    flow_id_mask = 65535;
    ring_id = 3;
    service_function = "sampling";
    sampling_rate = 10; //meaning 1:10 downsampling
    pipeline = 1;
    service_core_in_pipeline = 0;
    service_core_id = 2;
    kafka_topic = "topic_flow3";
},
```
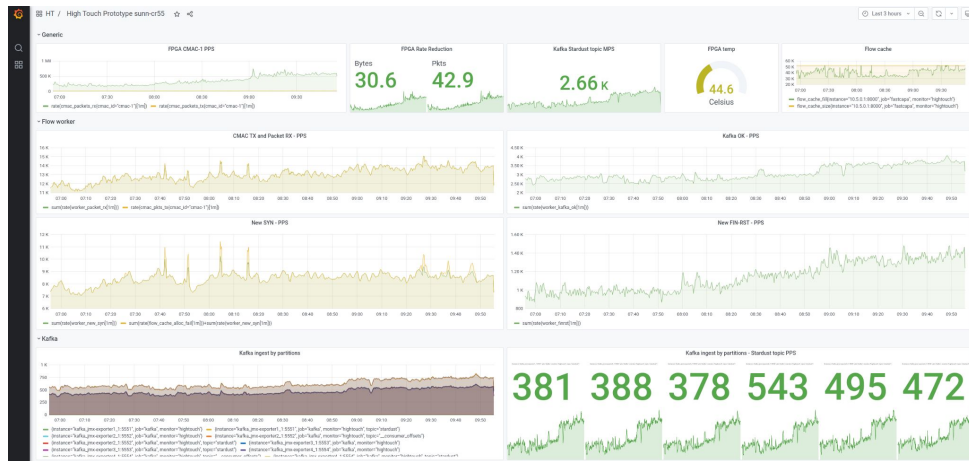
**Sampling**

```
{
    protocol = 6;
    protocol_mask = 255;
    srcIP = "192.168.25.4";
    srcIP_mask = 32;
    dstIP = "192.168.25.5";
    dstIP_mask = 32;
    srcPort = 10002;
    srcPort_mask = 10002;
    dstPort = 5202;
    dstPort_mask = 5202;
    priority = 102;

    flow_id = 2;
    flow_id_mask = 65535;
    ring_id = 2;
    service_function = "histogram";
    resolution = 100;//ns for inter arrival
    report_interval_tsc = 280000000; //CPU cycle count not actual
    //report_interval_tsc = 280; //CPU cycle count not actual
    pipeline = 0;
    service_core_in_pipeline = 1;
    service_core_id = 1;
    kafka_topic = "topic_flow2";
},
```

**Histogram**

```
{
    protocol = 6;
    protocol_mask = 0;
    srcIP = "0.0.0.0";
    srcIP_mask = 0;
    dstIP = "0.0.0.0";
    dstIP_mask = 0;
    srcPort = 0;
    srcPort_mask = 65535;
    dstPort = 0;
    dstPort_mask = 65535;
    priority = 1;

    flow_id = 0;
    flow_id_mask = 0;
    ring_id = 0; //drop at flow_worker
    service_function = "passthrough";///"drop"
    pipeline = 0;
    service_core_in_pipeline = 0;
    service_core_id = 0;
    kafka_topic = "topic_drop";
},
```
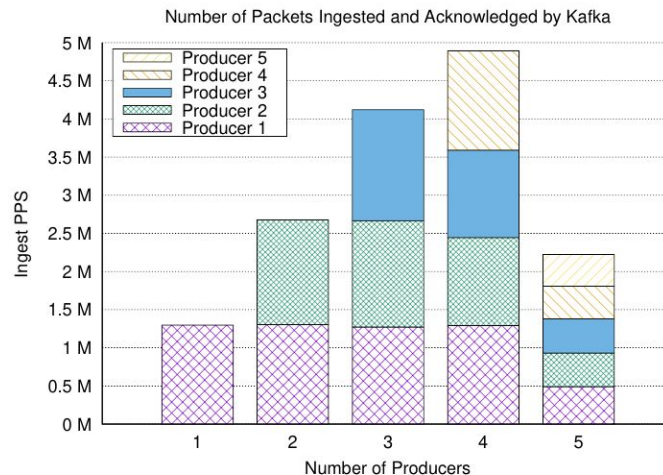
**Filter**

ESnet

# Fastcapa-ng Runtime Statistics



Fastcapa-ng pipeline statistics
Also in Grafana via Prometheus

# Kafka Performance

- *Apache Kafka*: open-source distributed stream platform.

- Docker-compose for a single server:

  – bitnami/kafka (x6), bitnami/zookeeper (x3)

  – bitnami/jmx-exporter

  – prom/prometheus

  – grafana/grafana

- 5M messages per second Kafka ingest performance

  demonstrated on single server.

- Possible bottlenecks to go higher:

  – Librdkafka C client (inside Fastcapa-ng)

  – Docker proxy - network
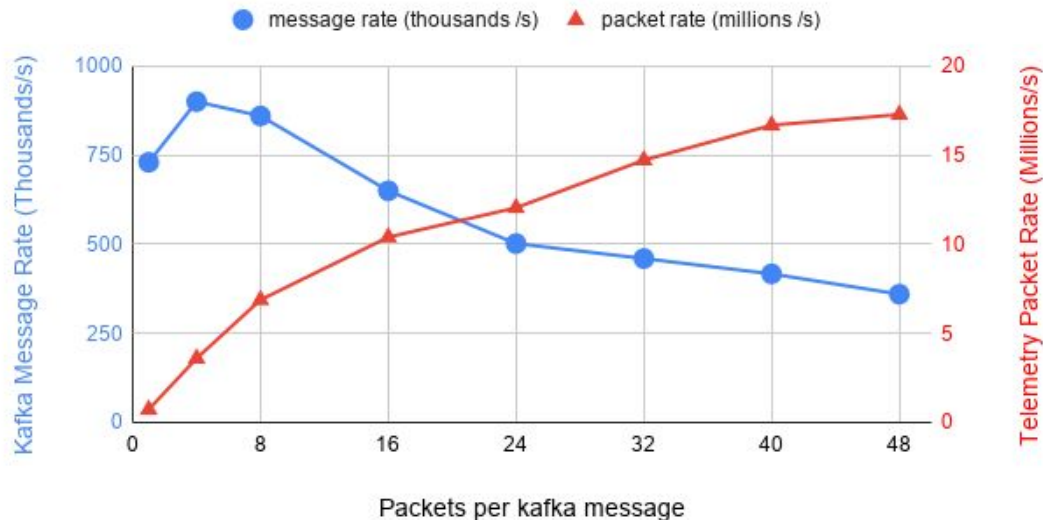
  – CPU - Client and brokers share the host



*~5M messages per second ingest*
*untuned single server / 6 broker / parallel producers*
*Kafka Benchmark tool, 64K message batches*

# Fastcapa's Kafka - going over 15M PPS

On top of **message batching** (handled by librdkafka), we need **packet batching** (handled by Fastcapa / client application).

*That means that one Kafka message contains multiple telemetry packets. Client application has to unpack.*



Single Kafka client performance using packet batching

ESnet

# High Touch Application Programming

- High Touch Applications can be implemented using **Kafka Streams** - an easy way to program real-time applications on stream of data.
- Expressive, highly scalable and fault tolerant API that allows: aggregation, filtering, counting, grouping data...

```
int THRES = 10;
KTable<Windowed<String>, Long> SYNcounts = stream
        .filter((k, telemetry) -> telemetry.isSYN())
        .groupBy((k, telemetry) -> telemetry.getIPDstAddr())
        .windowedBy(TimeWindows.of(Duration.ofSeconds(5)))
        .count(Materialized.with(String(), Long()))
        .filter((key, value) -> value > THRES);
SYNcounts.toStream().to("syn-attacks");
```

*Example: High Touch SYN Flood Detection*

# Conclusion

- We are processing millions of telemetry messages per second

- Data ingest is handled by **Fastcapa-ng**, an ESnet DPDK + Kafka project

  - Multi-stage, multi-pipeline architecture with easy configurability

  - Executes stateful functions: sampling, histogram creation, etc.

  - We can push 15M telemetry messages to Kafka with a single server

- Kafka streams: high-level application programming on telemetry streams

**We are working on open-sourcing fastcapa-ng: targeting Fall 2021.**

ESnet

# Questions…

richard@es.net

ESnet