# Flows and BGP, a multi-Terabyte story

Louis Poinsignon - Cloudflare

# Why monitoring

- Billing
  - Reducing costs

- Traffic engineering
  - Where should we peer?
  - Where should we set-up a new PoP?
  - Optimizing our network

- Anomaly detection
  - Troubleshooting
  - Proactive monitoring and predictions

# Ways to monitor

- SNMP

- Flows

- BGP

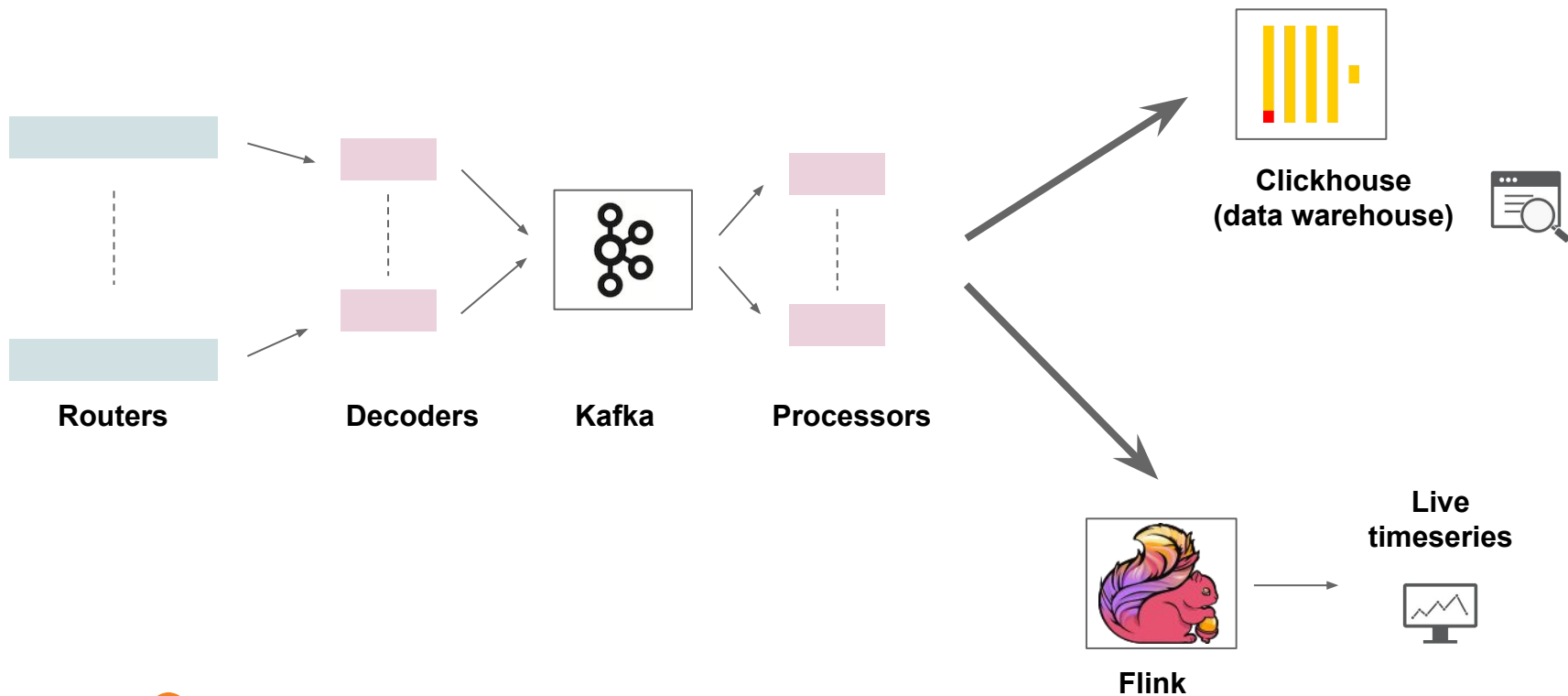What this presentation will be about

# Flows

# The previous pipeline

- Single machine
  - High CPU usage
  - Not easily accessible
- Dropping packets (around 10%)
- Not monitored
- sFlow cannot be aggregated with NetFlow
- Data not accessible to other teams / hard to develop on

- nfacctd + opentsdb

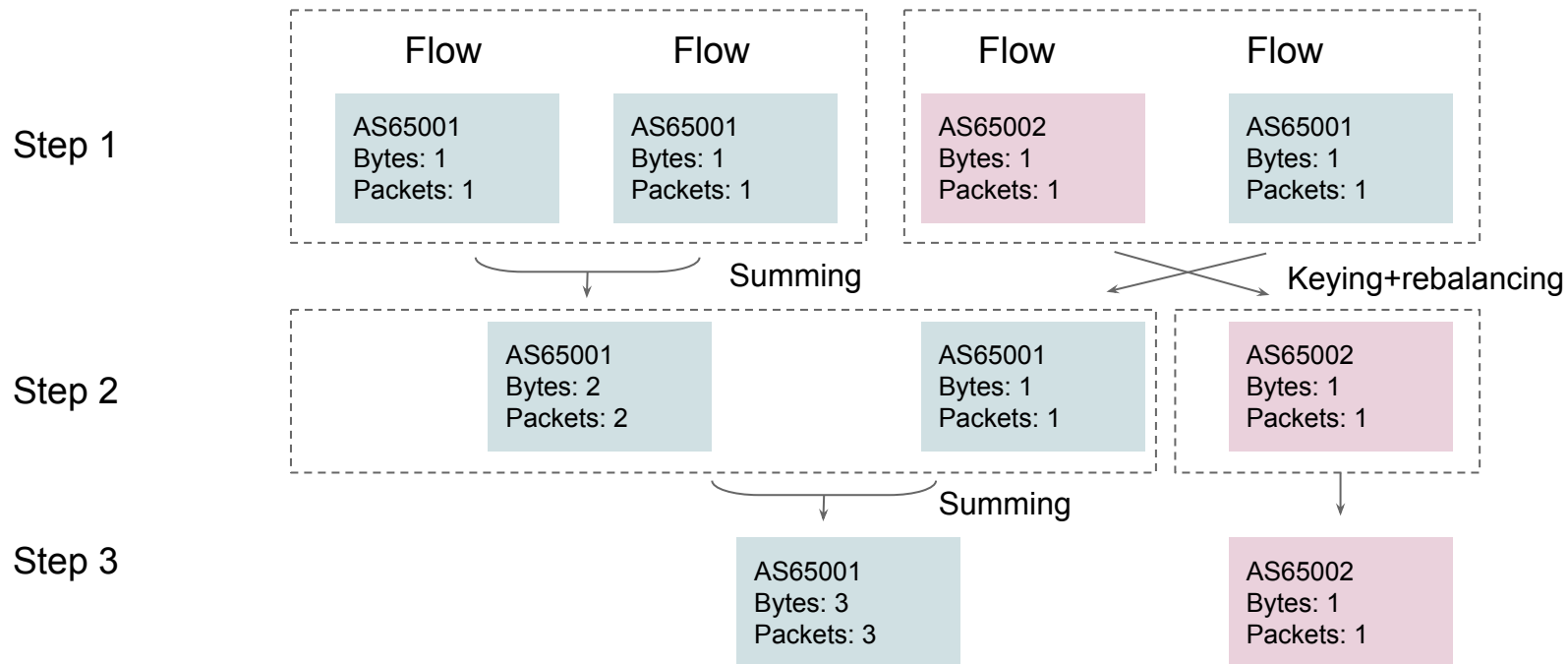CLOUDFLARE

# Time for an upgrade

- Let's use the tools provided by the platform team
    - Load-balancers
    - Mesos
    - Containers
    - Databases
    - Messages brokers
    - Big data processing clusters

CLOUDFLARE

# What we built



Routers          Decoders          Kafka          Processors

Clickhouse
(data warehouse)

Flink

Live
timeseries

CLOUDFLARE

# Flink - MapReduce

# Flink - Sample program

```
DataStream<FlowMessage> inData =
new FlinkKafkaConsumer09<FlowMessage>(
                "netflows-processed",
                new FlowMessageDeserializer(),
                propertiesConsumer);

DataStream<FlowMessage> inDataEyeball =
                inData.filter(new FlowFilter.EyeballFilter()).
                setParallelism(1).broadcast();

DataStream<FlowAggMessage> inDataAgg =
                inDataEyeball.map(new FlowUtils.Mapper("DstAS,colo"));

inDataAgg.reduce(new FlowTransformations.FlowAggReduceKey());
```
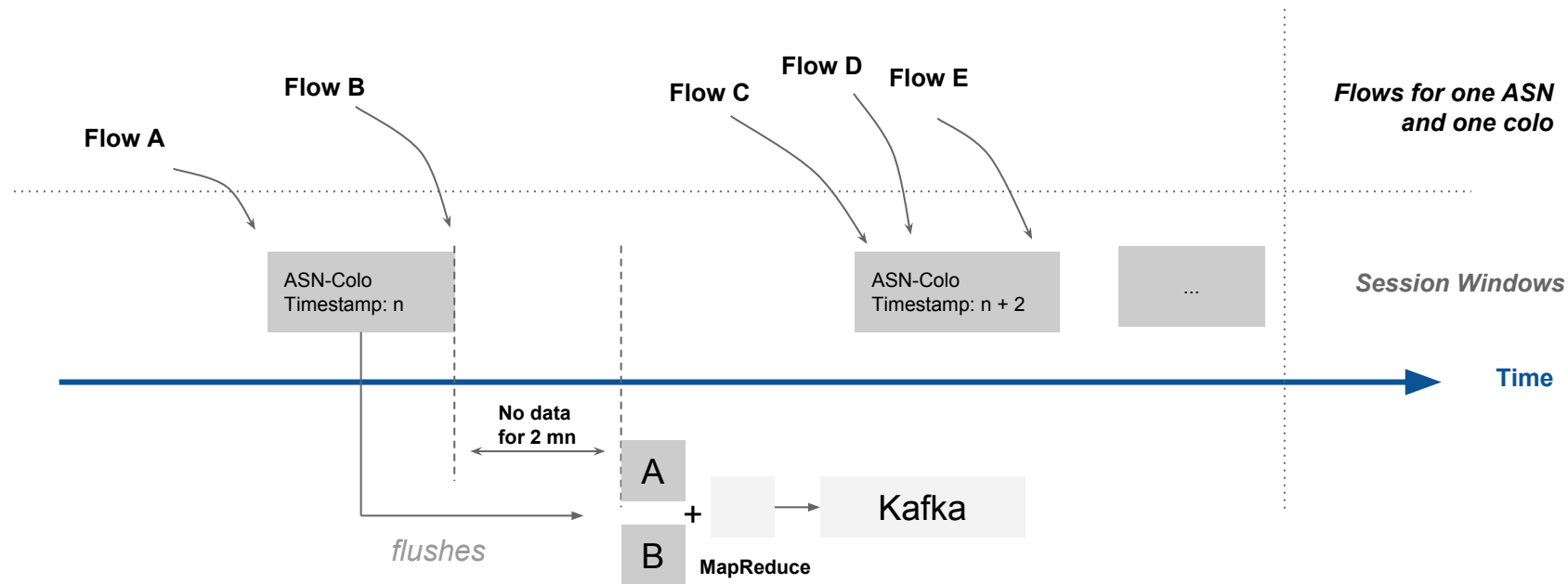
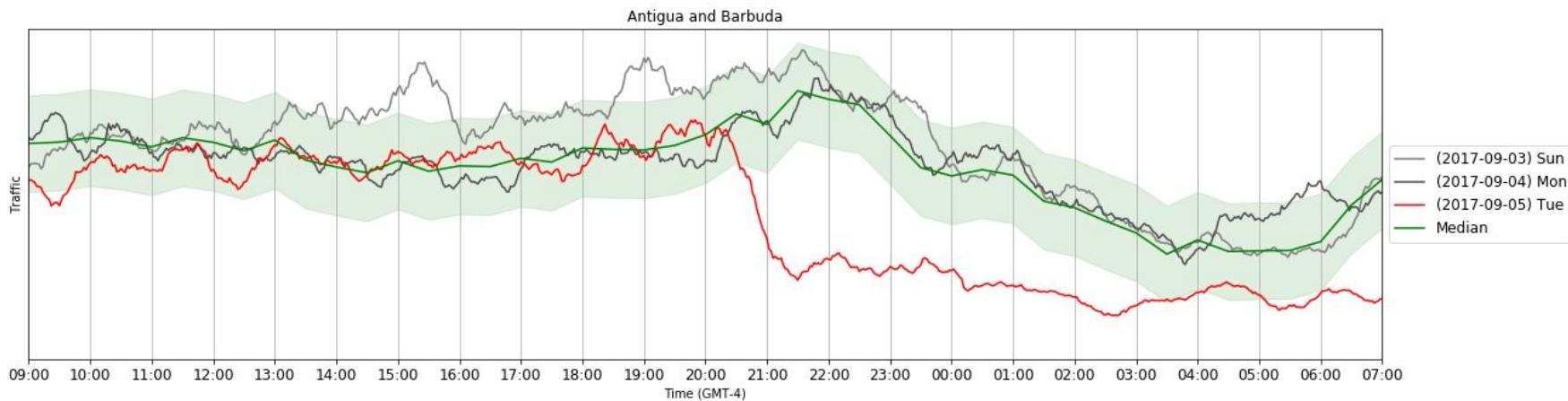Source (Kafka)

Filter

Mapping

Reduce



CLOUDFLARE

# Flink - Windowing

**Flow B**

**Flow A**

**Flow C**  **Flow D**  **Flow E**

*Flows for one ASN and one colo*

ASN-Colo
Timestamp: n

ASN-Colo
Timestamp: n + 2

...

*Session Windows*

**Time**

**No data for 2 mn**

A

+  →  Kafka

*flushes*

B  **MapReduce**

CLOUDFLARE

# Results - Aggregations

**Traffic generated by all the datacenters, by country, by network, etc.**



Antigua and Barbuda

Legend:
- (2017-09-03) Sun
- (2017-09-04) Mon
- (2017-09-05) Tue
- Median

Y-axis: Traffic
X-axis: Time (GMT-4)

09:00 10:00 11:00 12:00 13:00 14:00 15:00 16:00 17:00 18:00 19:00 20:00 21:00 22:00 23:00 00:00 01:00 02:00 03:00 04:00 05:00 06:00 07:00

https://blog.cloudflare.com/the-story-of-two-outages/

# Results - Statistics

Top providers per country, IPv6 penetration, etc.

| | Network | AS | Ratio IPv6 |
|---|---|---|---|
| 1 | Virgin Media | 5089 | 0% |
| 2 | BTnet (BT's UK IP Network - AS2856) | 2856 | 27.92% |
| 3 | Sky Broadband | 5607 | 69.77% |
| 4 | EE | 12576 | 5.19% |

CLOUDFLARE®

# Results - Example: maintenance

Building a list with the best hours for maintenance



Normalized traffic of our PoP over a day

# Algorithms

**Derivation**

**Correlation**

    **Pearson coefficient**

        Quantify the difference

**Median**

    Remove small variations

**Variance**

    Variation intensity



Outliers

Local variance following median
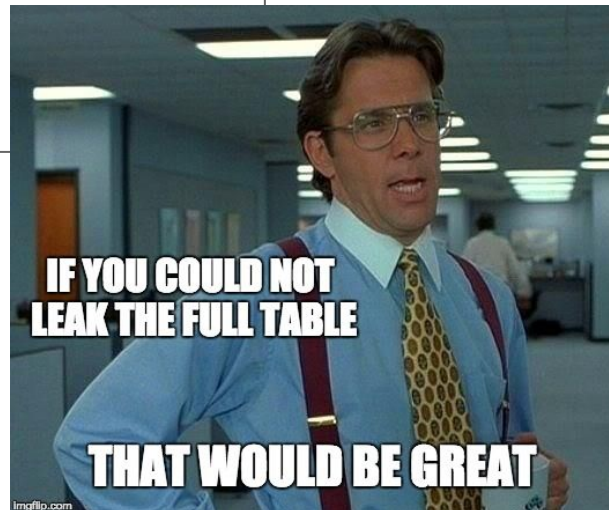
# BGP

# BGP flaps



Hi,

Your PeeringDB entry suggests 15000, but we have the maximum set to 20,000, hoe many prefixes are you sending?

Maybe we announce more than 20,000. At the moment I can not say for sure.
Please, increase more.

I can not say for sure.

Please, increase more.

IF YOU COULD NOT LEAK THE FULL TABLE

THAT WOULD BE GREAT

CLOUDFLARE

# BGP

100+ routers with full tables

Want to be aware of route-leaks or other anomalies

Need to do data analysis and periodic reports

Provide improved data for the flow pipeline
- Mappings prefixes to ASN
- Next-hop to ASN
- Peering/transit information

CLOUDFLARE®

# BGP Pipeline

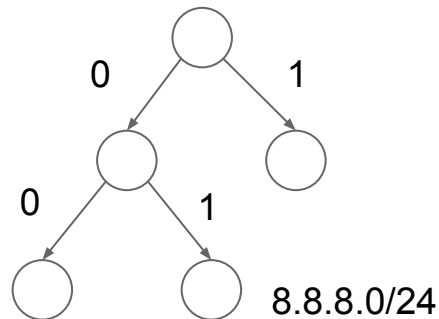Completely custom BGP implementation to live on a cloud with elastic IPs.

- More flexible than GoBGP
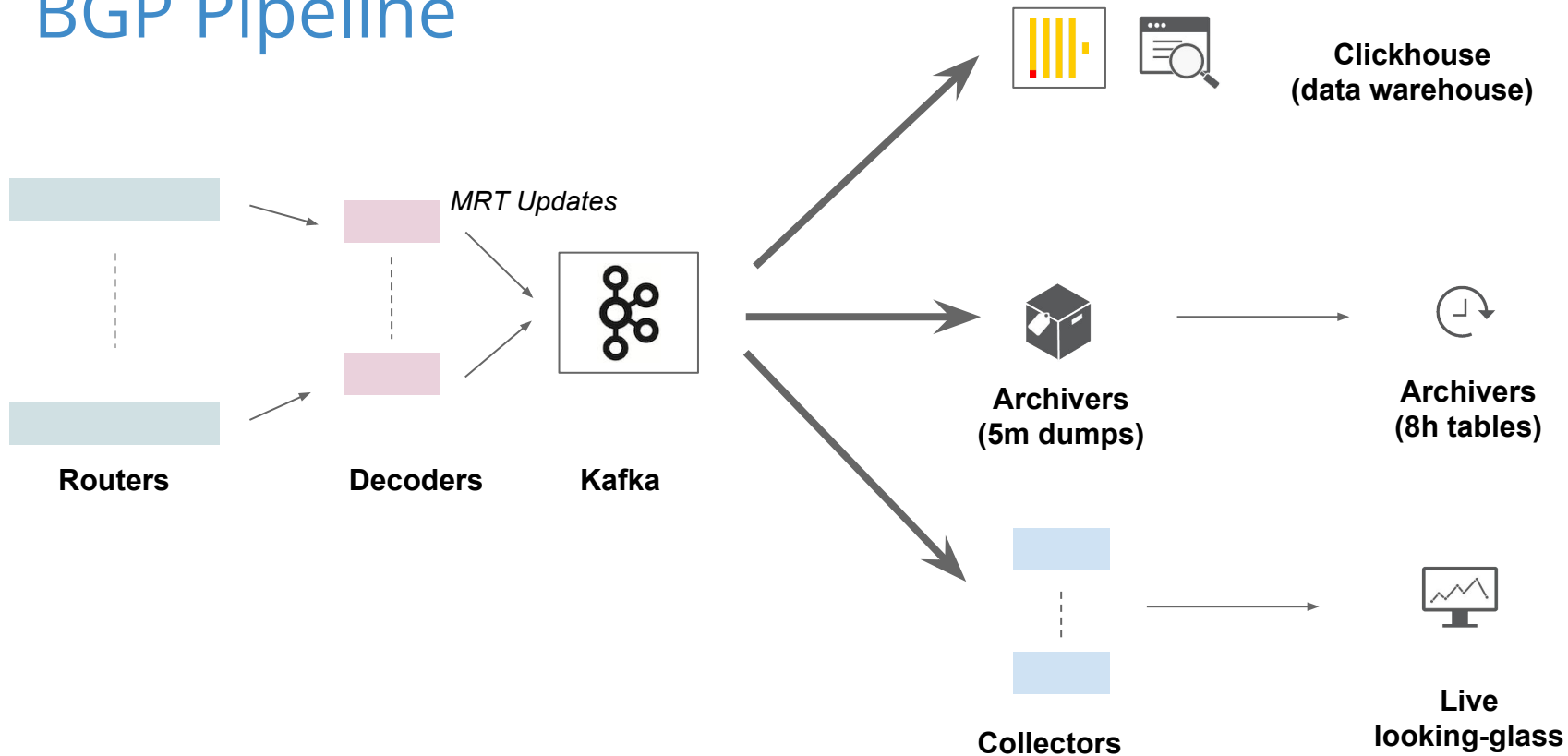
Horizontal scaling and hashing

- A single machine contains a subset of the data.
- Messages from x will always reach y

Develop custom trie implementation
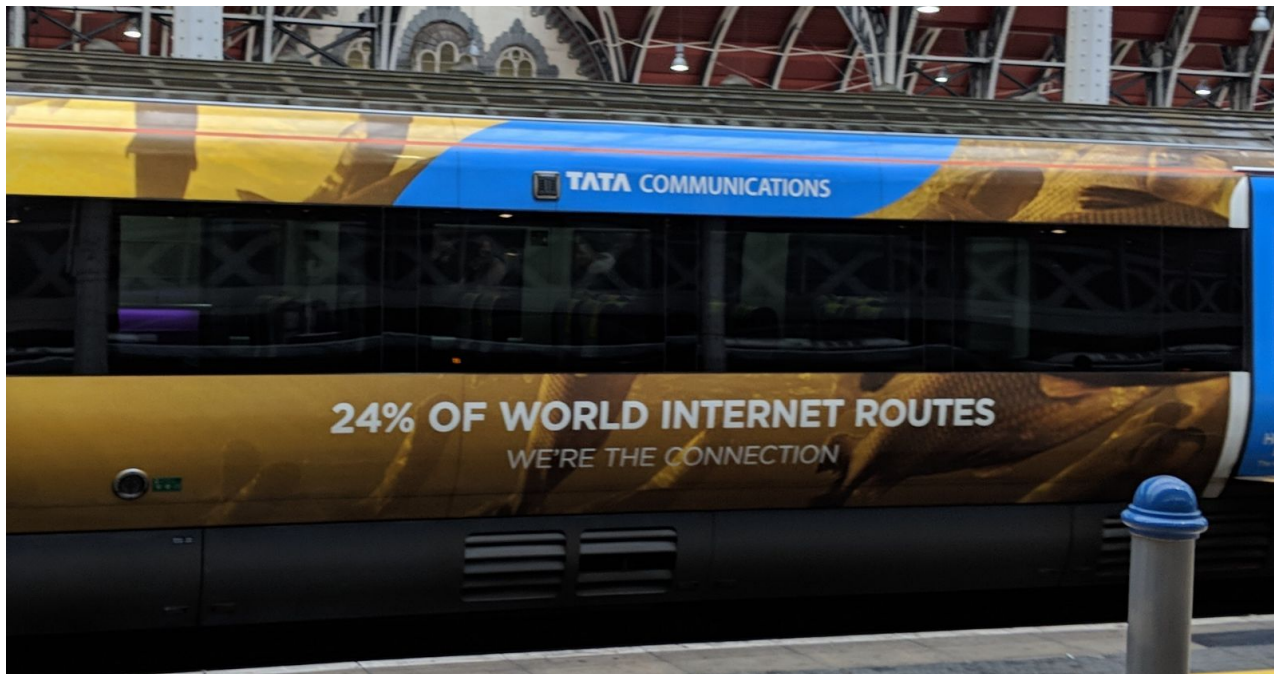
- Stores IP prefixes in an optimized way

0   1

0   1

8.8.8.0/24

# BGP Pipeline



*MRT Updates*

**Routers**

**Decoders**

**Kafka**

**Clickhouse
(data warehouse)**

**Archivers
(5m dumps)**

**Archivers
(8h tables)**

**Collectors**

**Live
looking-glass**

CLOUDFLARE

# BGP API

{
    "queried":
    "responded":
    "timeout": false,
    "ribs": {
        "        ": {
            "        ":
            "coloname": "EWR01",
            "routeraddr": "        ",
            "paths": [
                {
                    "prefix": "8.8.8.0/24",
                    "pathid": 0,
                    "nexthop": "        "
                    "origin": "IGP",
                    "med": 0,
                    "locpref": 200,
                    "communities": [
                        "        "
                    ],
                    "aspath": [
                        15169
                    ]
                }
            ]
        },

# Discoveries

- People sending us IX LAN prefixes

- Receiving smaller than /48 IPv6 and smaller than /24 IPv4

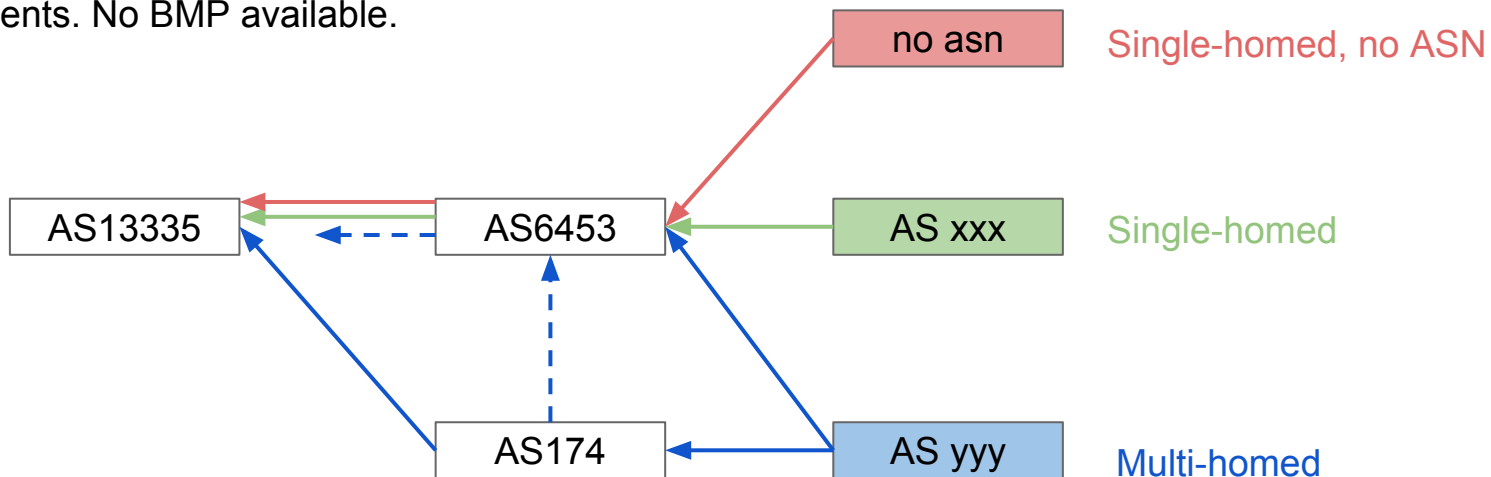- Longest AS-Path
    - 2402:8100:3980::/42 → 37 ASNs

# A train
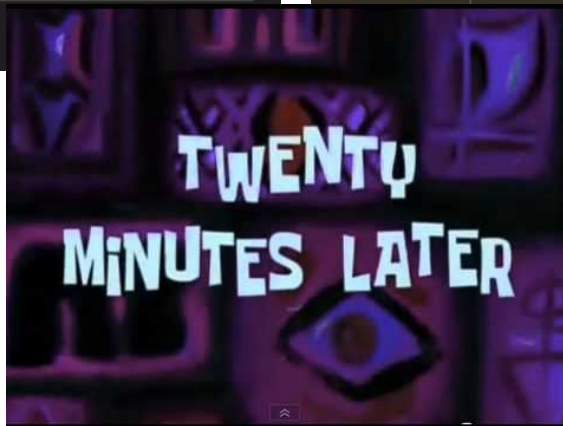


"We have approximately 170 000 routes"
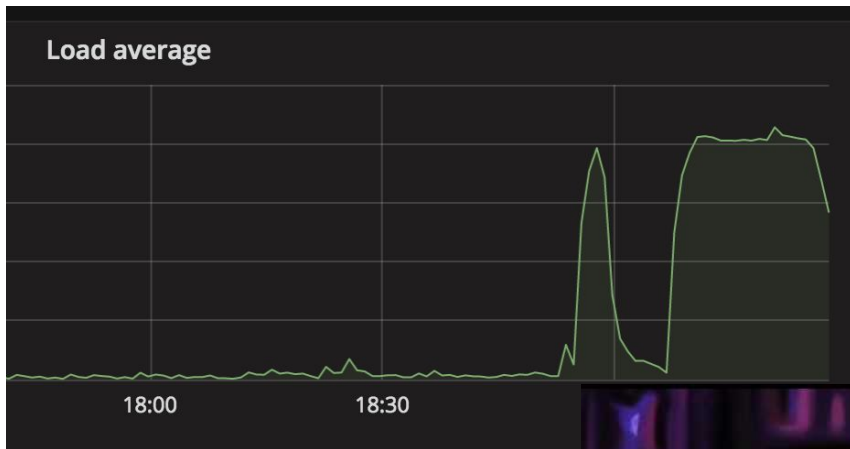
# TATA's Peeringdb

| | |
|---|---|
| **IPv4 Prefixes** | 220000 |
| **IPv6 Prefixes** | 11000 |

# Transit

3 types of clients. No BMP available.



no asn — Single-homed, no ASN

AS xxx — Single-homed

AS yyy — Multi-homed

AS13335   AS6453   AS174

CLOUDFLARE

# Let's process

# Results

**315997** routes (approx 42%)

- with 6453 as the last ASN with no transit listed behind
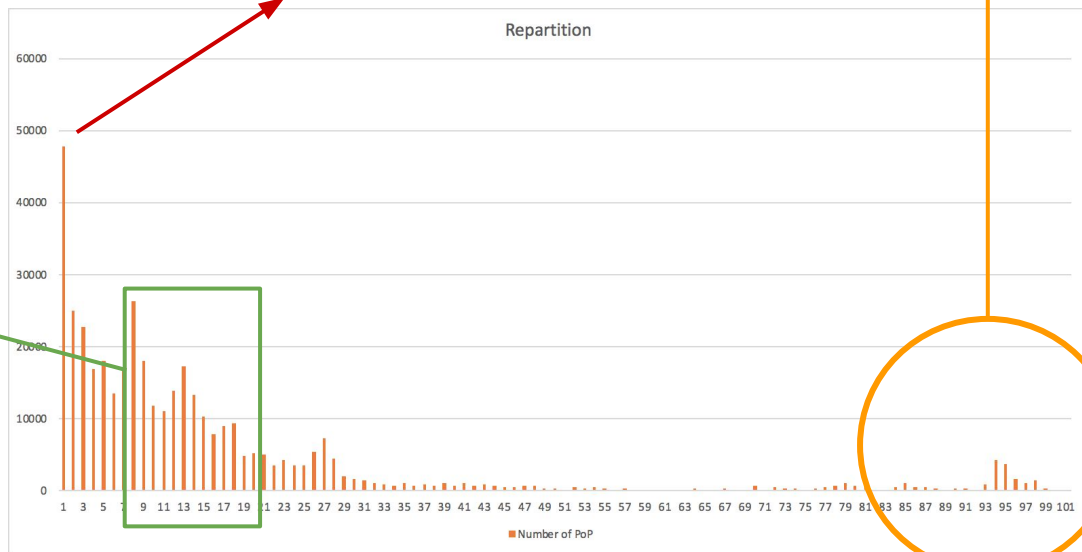
205 routes

- With 6453 as the last ASN
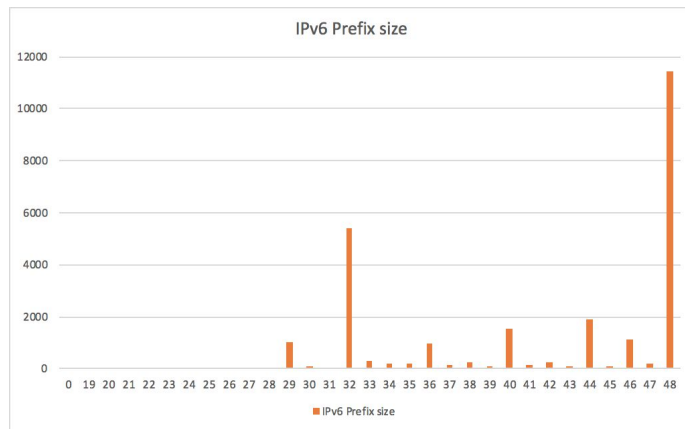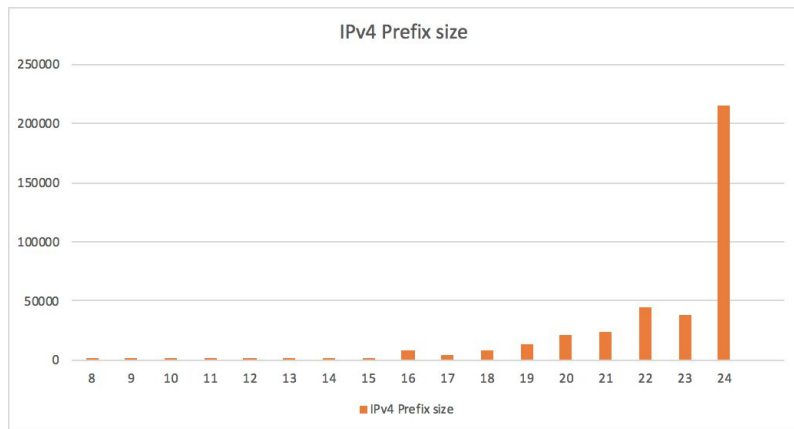
Around 100k routes per PoP

- With 6453 as ASN



Single-homed

Only as last resort

Around 24%

Repartition

# IP prefix sizes



IPv4 Prefix size



IPv6 Prefix size

# Size, costs

# Size

**Flows**

- Raw feed: 10-50 mbps
- Unique storage: 1-5 Terabytes
- Aggregated storage: 10-100 GigaBytes

**BGP**

- Raw feed: few kbps (100s updates per second)
- Around 35 GB of RAM used for storing the tables
- Around 100 MB per full table stored in MRT format

CLOUDFLARE

# Costs

The price of running it on a cloud (Amazon/Google/Azure) will depend on your size.

Rule of thumb:

| traffic | 10Gb/s | 100Gb/s | 1Tb/s |
|---------|--------|---------|-------|
| routers | 10 | 50 | 100 |
| Flows | $200/mo | $800/mo | $1,500/mo |
| BGP | $200/mo | $600/mo | $800/mo |

CLOUDFLARE®

Open-source

# Open-source (soon)

Flows:

- Decoders
- ~~Processors~~ → too specific to Cloudflare
- Flink example

BGP:

- Library to manage a session
- ~~Server, APIs~~ → too specific to Cloudflare

Tools:

- Maxmind DB encoder

Thank you!
Questions?