# New IP Who Dis?

Exploring AWS VPC Flow Logs with Immerse and OmniSciDB
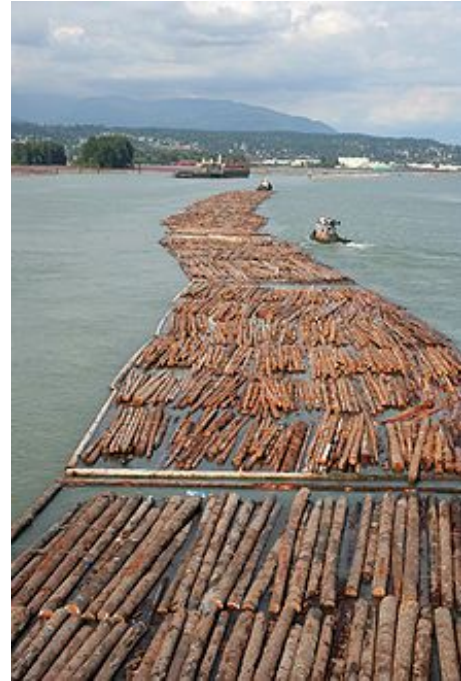
Tom Bowyer
Security Engineer- OmniSci

# Overview

- What is an AWS VPC Flow Log
- How do I turn it on
- What Information can I get from these logs
- QA

# What is an AWS VPC flow log

- Kinda like Netflow

- Contains all VPC network transit data

- Stores the data in S3 as a csv…..can also go to cloud watch if you hate yourself.

# How to turn on

# How to turn on

| account-id | Your AWS Account ID |
|---|---|
| interface-id | The ID of the network interface for which the traffic is recorded. |
| srcaddr | The source address |
| dstaddr | The Destination Address |
| srcport | The Source port |
| dstport | The Destination Port |
| protocol | What Protocol number |
| packets | Number of packets during the flow |
| bytes | Bytes Transferred during the flow |
| start | Start time of the flow (Unix) |
| end | End time of the flow (Unix) |
| action | Accepted or Rejected |



https://docs.aws.amazon.com/vpc/latest/userguide/flow-logs.html#flow-logs-fields

Additional metadata added just last month. Requires a format change in your flowlog config :( .

information from the log data.

When you create a new VPC Flow Log, in addition to existing fields, you can now choose to add the following meta-data:

- `vpc-id` : the ID of the VPC containing the source Elastic Network Interface (ENI).
- `subnet-id` : the ID of the subnet containing the source ENI.
- `instance-id` : the Amazon Elastic Compute Cloud (EC2) instance ID of the instance associated with the source interface. When the ENI is placed by AWS services (for example, AWS PrivateLink, NAT Gateway, Network Load Balancer etc) this field will be " `-` "
- `tcp-flags` : the bitmask for TCP Flags observed within the aggregation period. For example, `FIN` is 0x01 (1), `SYN` is 0x02 (2), `ACK` is 0x10 (16), `SYN` + `ACK` is 0x12 (18), etc. (the bits are specified in "Control Bits" section of RFC793 "Transmission Control Protocol Specification").
  This allows to understand who initiated or terminated the connection. TCP uses a three way handshake to establish a connection. The connecting machine sends a `SYN` packet to the destination, the destination replies with a `SYN + ACK` and, finally, the connecting machine sends an `ACK` . In the Flow Logs, the handshake is shown as two lines, with `tcp-flags` values of 2 ( `SYN` ), 18 ( `SYN + ACK` ). `ACK` is reported only when it is accompanied with SYN (otherwise it would be too much noise for you to filter out).
- `type` : the type of traffic : IPV4, IPV6 or Elastic Fabric Adapter.
- `pkt-srcaddr` : the packet-level IP address of the source. You typically use this field in conjunction with `srcaddr` to distinguish between the IP address of an intermediate layer through which traffic flows, such as a NAT gateway.
- `pkt-dstaddr` : the packet-level destination IP address, similar to the previous one, but for destination IP addresses.

# Now What

# ETL

- Add Geoenrichment
  - Need me some countries

- Add threat intel
  - OTX

- Fast import and fully automated

- Python very much

https://github.com/Zeerg/Conference-Talks/tree/master/omnisci-converge-2019/s3flow_sync

CONVERGE
by omni·sci

# Extract

Multiprocess step 1
- Pull data from S3 and insert into queue

```python
class S3Pull(Process):
    def __init__(self, **kwargs):
        super(S3Pull, self).__init__()
        self.sync_queue = kwargs.get('sync_queue')
        self.s3_client = boto3.client('s3')
        self.paginator = self.s3_client.get_paginator('list_objects')
        self.bucket = kwargs.get('bucket', None)
        self.bucket_prefix = kwargs.get('bucket_prefix', None)
        self.flow_date = kwargs.get('flow_date', None)

    def process_s3_files(self, bucket=None, key=None, date=None):
        logging.debug(f"Bucket: {bucket}")
        date_string = key + date
        logging.debug(f"Key to Iter: {date_string}")
        file_list = self.paginator.paginate(
                                    Bucket=bucket,
                                    Prefix=date_string,
                                    PaginationConfig={'MaxItems': 20000}
                                    )
        keys = []
        for page in file_list:
            for key in page['Contents']:
                keys.append(key['Key'])
        key_count = len(keys)
        count = 1
        for flow_log in keys:
            obj = self.s3_client.get_object(Bucket=bucket, Key=flow_log)
            df = pd.read_csv(io.BytesIO(obj['Body'].read()), compression='gzip')
            self.sync_queue.put(df)
            logging.info(f"Putting frame {count} of {key_count} into queue")
            count += 1
            time.sleep(.5)          You, 4 days ago • working poc

    def run(self):
        logging.info("Starting S3 Sync Process")
        self.process_s3_files(self.bucket, self.bucket_prefix, self.flow_date)
```

# Transform(ers)

Multiprocess step 2
- Pull from queue then add dem countries using maxmind
- Also mask things we don't want to share

- Add threat intel lookup

- Add to queue

```python
class PandaTransform(Process):
    def __init__(self, **kwargs):
        super(PandaTransform, self).__init__()
        self.mask_values = kwargs.get('mask', [])
        self.sync_queue = kwargs.get('sync_queue')
        self.transform_queue = kwargs.get('transform_queue')
        self.invalid_chars = ['-']
        if kwargs.get("mmdb", None):
            self.mmdb_geo = geoip2.database.Reader(kwargs.get("mmdb"))

    def mmdb_lookup(self, src_ip):
        try:
            response = self.mmdb_geo.city(src_ip)
            return (
                float(response.location.longitude),
                float(response.location.latitude),
                str(response.city.name),
                str(response.subdivisions.most_specific.name),
                str(response.postal.code),
                str(response.country.name),
                str(response.country.iso_code),
            )
        except Exception as e:
            logging.error(f"Failed to lookup {e}")
            return None

    def transform_files(self):
        while True:
            item = self.sync_queue.get()
            df = pd.DataFrame(item)
            df.replace({r: "xxxxxxxxxx" for r in self.mask_values}, regex=True, inplace=True)
            big_tup = list(df.itertuples(index=False, name=None))
            geo = False
            for log_tup in big_tup:
                log_list = log_tup[0].split(" ")
                log_struct = LogStruct()
                ip_address = log_list[3]

                # Build our tuple
                log_struct version = log list[0]
```

# Load

Multiprocess step 3
- Grab from load queue and...
- Load them frames



```python
import pandas as pd
import time

from multiprocessing import Process, JoinableQueue


class OmnisciLoader(Process):
    def __init__(self, **kwargs):
        super(OmnisciLoader, self).__init__()
        self.transform_queue = kwargs.get('transform_queue')
        self.table_name = kwargs.get('table_name')
        self.db_connection = kwargs.get('omnisci_connection')
        self.batch_size = 1000        You, 4 days ago • working poc

    def insert_data(self):
        log_list = []
        while True:
            frame_tuple = self.transform_queue.get()
            log_list.append(frame_tuple)
            self.transform_queue.task_done()
            if len(log_list) >= self.batch_size:
                df = pd.DataFrame(log_list)
                logging.info("Loading Flow Log Batch Into Table")
                try:
                    self.db_connection.load_table_columnar(self.table_name, df, preserve_index=False)
                    log_list = []
                except Exception as e:
                    logging.error(f"Fail to insert data {e}")
                    log_list = []
            if self.transform_queue.empty():
                logging.info("Queue Empty")
                time.sleep(.8)

    def run(self):
        logging.info("Starting the Omniscidb Load Process")
        self.insert_data()
```

# ETL



https://github.com/Zeerg/Conference-Talks/tree/master/omnisci-converge-2019/s3flow_sync

# What about the ASN?

- ?????? uhhhhh
- ETL 150 million records again adding 1 column

- Or

- Do a join on a table that has ASNs and IPs



WHEN YOU DON'T HAVE TIME TO ETL 150MILLION ROWS

Improvise. Adapt. Overcome

imgflip.com

CONVERGE
by omni·sci

# What about the ASN?

# What about the ASN?

```
omnisql> \d converge_flowlogs
Table converge_flowlogs doesn't exist
omnisql> \d converge_flowlog
CREATE VIEW converge_flowlog AS SELECT *
FROM flowlogs
LEFT JOIN uniq_with_asn ON flowlogs.source_address=uniq_with_asn.column_1;

View columns:

version INTEGER,
account_id TEXT ENCODING DICT(32),
interface_id TEXT ENCODING DICT(32),
source_address TEXT ENCODING DICT(32),
dest_address TEXT ENCODING DICT(32),
src_port TEXT ENCODING DICT(32),
dest_port TEXT ENCODING DICT(32),
protocol TEXT ENCODING DICT(32),
packets INTEGER,
bytes INTEGER,
flowlog_start TIMESTAMP(0),
flowlog_end TIMESTAMP(0),
action TEXT ENCODING DICT(32),
log_status TEXT ENCODING DICT(32),
src_lon FLOAT,
src_lat FLOAT,
src_city TEXT ENCODING DICT(32),
src_state TEXT ENCODING DICT(32),
src_zip_code TEXT ENCODING DICT(32),
src_country TEXT ENCODING DICT(32),
src_country_iso TEXT ENCODING DICT(32),
otx_intel BOOLEAN,
rowid BIGINT NOT NULL,
ip_int BIGINT,
column_1 TEXT ENCODING DICT(32),
network_int BIGINT,
broadcast_int BIGINT,
ip_as_int BIGINT,
asn TEXT ENCODING DICT(32),
country TEXT ENCODING DICT(32),
org TEXT ENCODING DICT(32),
rownotusd INTEGER,
rowid0 BIGINT
omnisql>
```
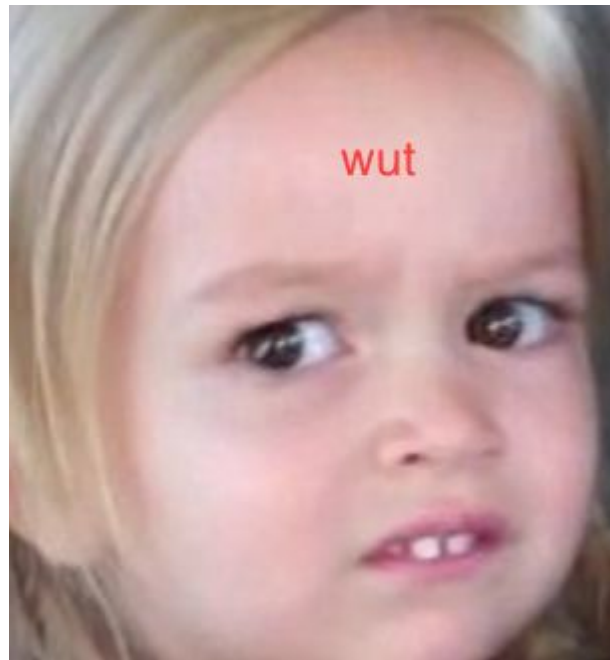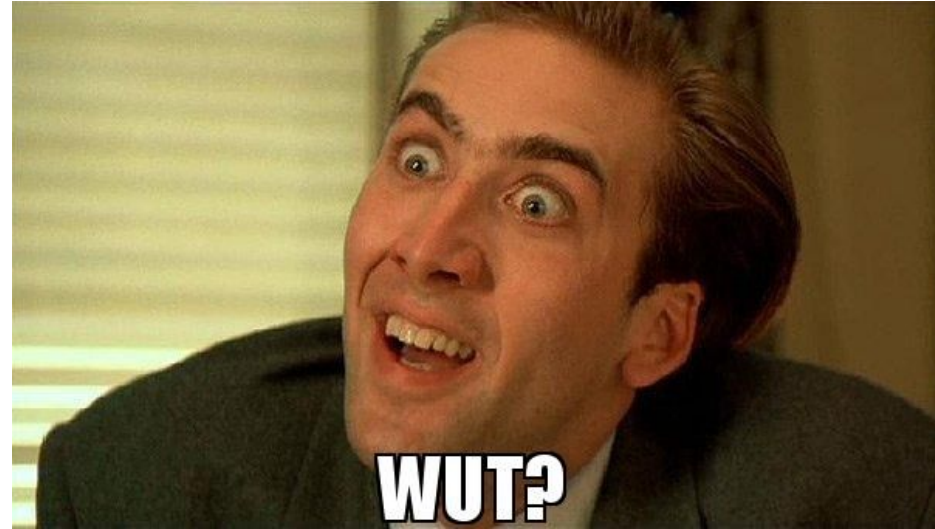
# Now We Explore

# Improving Security Posture

- Where to begin?

- Hey we need to limit outbound ports for this awesome compliance requirement.

- or really removing 0.0.0.0/0 outbound



CONVERGE
by omni·sci

# Threat Hunting



- Checking for allowed SSH

- Are we using port 80 inbound still?

- Is MySQL open to the world

- Threat Hunting

- Do I have hosts that allow ping

# Pro Tip

- If network ACLs attached to a NAT gateway don't explicitly deny traffic from the internet, internet traffic to the NAT gateway appears accepted.

- However, the actual traffic isn't accepted by the NAT gateway and is dropped.

- Heart Attack No More!

# Links

- ASN List

https://iptoasn.com/

- Maxmind

https://www.maxmind.com/en/geoip2-city

- ETL Tool

https://github.com/Zeerg/Conference-Talks/tree/master/omnisci-converge-2019/s3flow_sync

Q&A

thank you