

# Cats v Dogs

Data Streaming on AWS

**Alice Moyon**

[github.com/alicemoyon](https://github.com/alicemoyon)  
[linkedin.com/in/alicemoyon](https://linkedin.com/in/alicemoyon)  
[alicemoyon.ie](https://alicemoyon.ie)



Photo credit: [kevin turcios](#) on [Unsplash](#)

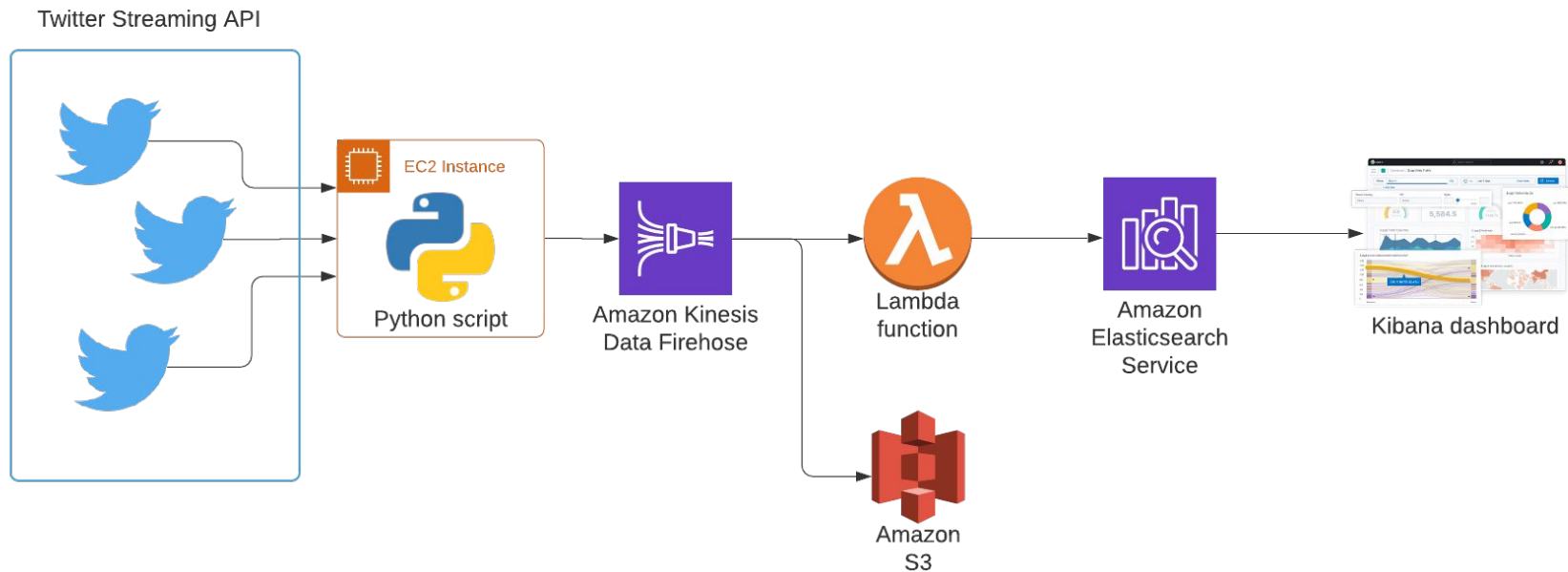
# Project objectives

- Learn the fundamental components of a data streaming pipeline
- Find a simple, highly abstracted cloud-based method to analyse a data stream in real-time
- Visualise real-time data on a dashboard

# Application objectives

- Stream tweets containing image or videos of cats and dogs
- Analyse tweets in real-time to determine most popular of the two
- Visualise evolution on a dashboard updated in real-time

# Architecture



# The new Twitter API v2

- Wide ranging and very well documented - a world of possibilities!
- Official SDKs for Python, Node.js, Ruby + community libraries
- Two streaming options:
  - Sample Stream - random 1% sample of all tweets
  - Filtered Stream - filter with rules
- Returns Tweet objects in JSON
- Free, with some limitations:
  - Rate limits
    - 500,000 Tweets per month
    - 25 concurrent rules
    - 50 connections requests per 15-min window
    - 450 filter adds or deletes per 15-min window
  - Restricted features
    - Redundant connection & backfill
    - Advanced operators for rules (place, bio...)

```
{
  "data": {
    "id": "1417910215564046336",
    "public_metrics": {
      "retweet_count": 1,
      "reply_count": 0,
      "like_count": 0,
      "quote_count": 0
    },
    "text": "RT @WendyDranfield: Joey's keeping cool in the long grass. 🌞🐾💕 #Cats #heatwave https://t.co/LubJdTIPKM"
  },
  "matching_rules": [
    {
      "id": "1417558067625988000",
      "tag": "cats with images"
    }
  ]
}
```

# The producer script - Python on EC2

- Sets up the filtering rules for the Twitter stream
- Opens a persistent HTTP connection to the API
  - Query parameters correspond to the tweet fields we want to retrieve.
- Puts the records into an AWS Firehose delivery stream as they come in through the API connection
- Issues encountered:
  - Trial & error with the query parameter syntax and filtering rules
  - Twitter bug: 0 bytes record every 20 minutes, requiring manual interrupt and reconnection
    - Wrapped in a loop with automatic reconnection if any exceptions occur
- Developed and tested on local machine then deployed to an AWS EC2 instance (Ubuntu).

# Firehose + Lambda for in-stream processing

- Kinesis Data Firehose:
  - AWS Service that can ingest, transform and deliver streaming data in real-time to S3, Redshift, Elastic Search, HTTP endpoints + others
    - Serverless, fully managed
    - Automatically scales to match data throughput
    - Loads to destination within 60s
  - Very easy to test in the AWS console with default record provided
  - Associate with a Lambda function for in-stream processing before delivery to destination
- Lambda:
  - Transform and process data records
  - Must include decoding & encoding again
  - Outputs back to the Firehose stream for delivery to destination
- Delivery to AWS Elasticsearch Service + S3 backup
  - Issue: records are saved to S3 in their original state
- Monitoring: enable AWS CloudWatch + see failed records on S3

# More on Lambda function

- Use Python, Java, Node.js, Go etc...
- Very versatile: can be used for lots of scenarios, triggered by specific events
- Create *Layers* to add Python packages used in the function
- Decode/Encode:
  - Records are Base64 encoded as they are ingested into the Firehose stream ([bytes-to-text encoding](#))
  - Decode before processing, re-encode before returning the processed record to the stream



lambda\_function ×

Execution results ×



```
1 import base64
2 import json
3
4
5 def lambda_handler(event, context):
6     output = []
7     for record in event['records']:
8         # decode the record data
9         payload = base64.b64decode(record['data']).decode()
10        # convert decoded string into json
11        payload = json.loads(payload)
12        print(type(payload))
13        # clean up the record
14        cleaned_payload = dict(tweet_id=payload['data']['id'],
15                               retweets=payload['data']['public_metrics']['retweet_count'],
16                               replies=payload['data']['public_metrics']['reply_count'],
17                               likes=payload['data']['public_metrics']['like_count'],
18                               category=payload['matching_rules'][0]['tag'],
19                               timestamp=record['approximateArrivalTimestamp']
20        )
21        # convert dictionary to json string and re-encode the payload
22        new_payload = base64.b64encode(json.dumps(cleaned_payload).encode())
23        output.append({
24            'recordId': record['recordId'],
25            'result': 'Ok',
26            'data': new_payload
27        })
28    return {'records': output}
29
```



# Testing the Lambda function

## Configure test event

A function can have up to 10 test events. The events are persisted so you can switch to another computer or web browser and test your function with the same events.

- ☐ Create new test event
- ☒ Edit saved test events

### Saved Test Events

testEventWithDoubleQuotes

```
1 {  
2   "invocationId": "invocationIdExample",  
3   "deliveryStreamArn": "arn:aws:kinesis:EXAMPLE",  
4   "region": "eu-west-1",  
5   "records": [  
6     {  
7       "recordId": "49546986683135544286507457936321625675700192471156785154",  
8       "approximateArrivalTimestamp": 1495072949453,  
9       "data": "eyJkYXRhIjoeyJpZCI6ICIxNDA2MjY1MDA4OTc2MDU2MzIwIiwgInB1YmtpY19tZXRYaWNW":  
10    }  
11 ]  
12 }  
13
```

Encode your data record separately, then use the encoded string in your test event

# Another option for in-stream analytics - Kinesis Data Analytics

## Kinesis Data Analytics vs transformation Lambdas

### Transformational Lambda

- Python + Pandas
- Filter / aggregate
- Fixed window
- Great for data transformations per item
- Cannot combine multiple streams
- Not the best way to send output to another destination

### Kinesis Data Analytics

- SQL
- Filter / aggregate
- We control the window
- Lets us look at the stream in chunks
- Can combine multiple streams
- Can send output to another stream or other destinations

*Source: DataCamp (see references slide)*

# AWS Elasticsearch Service



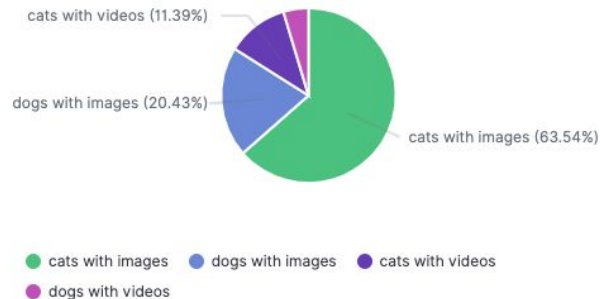
- Seamless implementation, High level of abstraction with AWS ES
- So I know very little...but in short
  - AWS Service is fully managed to deploy and run an Elasticsearch cluster
  - Elasticsearch is a distributed, JSON-based search & analytics engine
  - Built-in *Kibana* for visualisation & near real-time dashboarding
  - Easy set-up through AWS Console, automatic *indexing* with Firehose
  - AWS version < full capabilities
- Not super intuitive at first but looks like it has amazing capabilities if you can dig deeper and explore
- Issues encountered:
  - Needs strict JSON format
  - Data type change in the index
  - No fine-grained access control with Free Tier ES

# Kibana dashboard

- Very easy to create dashboard if only using basic aggregation
  - (once you figure out the jargon)
  - Index patterns (i.e. schemas) pulled from Elasticsearch
  - First, create individual visualisations
    - Great choice out of the box incl. geovisualisations
  - Then add visualisations to a dashboard
- Uses Elasticsearch index patterns
  - Automatically pulled from Elasticsearch
  - Watchout: data types cannot be edited
- For more, need to learn specific scripting language designed for Elasticsearch (Painless, Lucene and JSON style requests)
  - Can create scripted fields for visualisations
  - Myriads of possibilities for customisation

# So who's winning?

## So Far Today ①



## Total Tweets

7,709

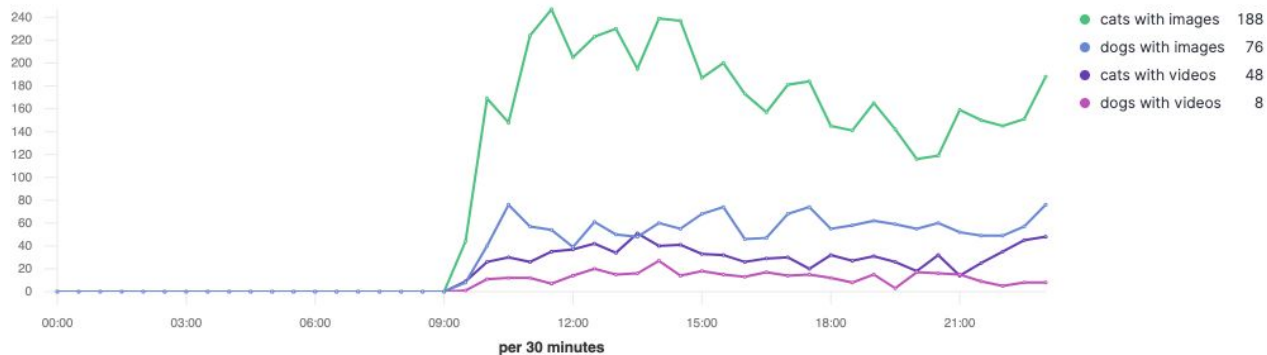
## Total Dog Tweets ①

1,933

## Total Cat Tweets ①

5,776

## Tweet timeline today



# Estimated AWS cost for this project

- Firehose
  - Charged only for the data transmitted
  - $N \text{ records} * \text{record size (round up to nearest 5KB)}$
  - \$0.031 per GB for first 500TB/month
  - Here:  $5\text{KB} * 450\text{k records} = \text{approx. } 2.25\text{GB/mth} \rightarrow \text{less than \$1}$
- Lambda
  - Charged based on number of requests and code execution duration
  - Need  $< 450\text{k requests/mth}$ , Free Tier includes  $1\text{M requests/mth} \rightarrow \text{Free}$
- AWS Elasticsearch service
  - Using Free Tier (instance type `t2.small.elasticsearch`)  $\rightarrow \text{Free}$
- EC2
  - Free Tier `t2.micro` for 12 months  $\rightarrow \text{Free}$
- S3
  - $450\text{k PUT requests/month} \rightarrow \text{\$2.27}$
  - Storage  $< 5\text{GB} \rightarrow \text{Free}$

# That was Phase 1.....

## Possible future enhancements

- Including number of likes in the analysis
  - Loop to continually search for tweets already streamed and update their like count
- More complex processing
  - e.g. exclude promotional tweets, include more fields, sentiment analysis
- Additional data source
  - E.g. Reddit
  - In-stream image analysis with call to Computer Vision model
- From Elastic to a web page
- Explore Elastic & Kibana capabilities
  - More complex analysis in Elastic
  - More sophisticated dashboard visualisations
  - Conditional text based on the real-time info? Is that possible?

# Questions / Discussion

- Has anybody used Elasticsearch or Kibana before?
- Or different tools for a similar use case?
- Any suggestions for improvements?
- What frameworks can be used to implement real-time updates in a web app?



# Some useful references

- Twitter filtered stream API  
<https://developer.twitter.com/en/docs/twitter-api/tweets/filtered-stream/introduction>
- DataCamp tutorial  
<https://learn.datacamp.com/courses/streaming-data-with-aws-kinesis-and-lambda>
- Base64 Wikipedia definition <https://en.wikipedia.org/wiki/Base64>