

Real-Time Data Streaming with Python + AWS Kinesis — How To... (part 1)



Tom Thornton

[Follow](#)

May 22, 2019 · 14 min read

In this article we will be focusing on the use of AWS Kinesis with Python and Node.js to stream data in near real-time to ElasticSearch.

Kinesis provides the infrastructure for high-throughput data processing and analytics. We can leverage this to push data through to storage/visualisation services for aggregation and historical interrogation.

In this guide we will be using Python 3.6 and AWS' boto3, pandas and inbuilt functions.

In Part 1, we will discuss each of the segments of the Kinesis service, what you can use them for and finally walkthrough a worked example of streaming CSV data with Kinesis to AWS' ElasticSearch service for visualisation in Kibana dashboards. This will cover the necessary code to stream data from source to a Kinesis Data Stream and to transform/push data to ElasticSearch using a Kinesis Firehose.

Every AWS service I will be working with are all within the same region (eu-west-1 for me) so you may have to make adjustments to suit your purposes.

. . .

Kinesis Data Streams

Data streams are like endpoints which accept batches of data to be processed. They provide “shards” which you can think of like processing nodes. The more shards, the more records you can process in parallel and therefore reducing total runtime to process the same number of records. The Data Stream ultimately provides a quick and easy way of bussing your data to a consumer. Too little shards and your data may end up in a dead letter queue, which means you don’t have enough shards for the volume of

data and they are becoming saturated thus rejecting any further data. By default Data Streams will hold data for 24 hours, but this can be modified up to 168 hours. This allows your applications to continue processing items in the Data Stream should the consumer have an issue. You can also provide a timestamp when querying for data from the stream to play back events that have been processed in the past.

• • •

Kinesis Firehoses

Firehoses are the opposite of Data Streams, they distribute records/data to specified endpoints configured in the Firehose from a specified source. In a Firehose, we can also specify any transformations or formatting we'd like to apply to the data before it's pushed out. We will be using Firehose later on to connect our Data Stream to ElasticSearch and provide the transformation to JSON logic.

Firehose provides the ElasticSearch Service domain as an output, but you can also target database services or just use a HTTP request. You would be using HTTP requests to submit data to a non-AWS ElasticSearch domain.

• • •

Getting Started...

Creating a Data Stream is easy, all you have to do is provide a name for the stream and a number of shards to use behind the stream. You can use the handy estimator to figure out how many you will need.

To do this simply, take the average size of the source data you'll be posting and divide this by the number of records to get your average record size. Once you have this you can feed it into the estimator to get a guide for the number of shards you'll need. I'd suggest slightly over-estimating to avoid any potential hiccups:

▼ Estimate the number of shards you'll need

Shard calculator

Average record size	<input type="text" value="2"/> KB
Record size is an integer between 1 and 1024	
Max records written	<input type="text" value="1000"/> per second
(Number of records per second) x (Number of producers)	
Number of consumer applications	<input type="text" value="1"/>

Estimated shards [Use this value](#)

Feed the estimated shards value into the Number of shards, you'll see a single shard can provide 1MB of writes/second, and 2MB of reads/second. This scales linearly as you increase shards.

Number of shards*

You can provision up to 500 more shards before hitting your account limit of 500. [Learn more or request a shard limit increase for this account](#)

Total stream capacity Values are calculated based on the number of shards entered above.

Write MB per second

1000 Records per second

Read MB per second

Using this approach, we can clearly see that Kinesis Data Streams are much more efficient when working with smaller amounts of data at lower throughput. We will be looking to use Kinesis for high-throughput processing, so we must try to reduce our data payload as much as possible to keep the costs low. More shards = More Cost.

Full details on Kinesis pricing can be found here —
<https://aws.amazon.com/kinesis/data-streams/pricing/>

After hitting create, it'll take a few moments to create the Data Stream, but once it's done you should see a status of 'Active':

Stream DataStream has been created

To enable server encryption, add shard-level metrics, or modify the shard number, view the [details](#).

[Create Kinesis stream](#) [Connect Kinesis consumers](#) [Actions](#)

Filter Kinesis streams

	Kinesis stream name	Number of shards	Status	Created
<input type="checkbox"/>	DataStream	1	Active	0

We're now ready to start pushing data to the Data Stream!

Creating an ElasticSearch Service Domain (optional)

For this guide, we will be using an ElasticSearch Service domain provided by AWS. You can of course use an existing ElasticSearch instance you have, but for simplicity we will be sticking to the services provided by AWS. I'll guide through setting up a simple ElasticSearch domain, as we will be using this later with a Kinesis Firehose.

Go to the ElasticSearch Service, and hit 'Create a new domain'...

Here I would just select 'Development and testing' and for the version, select the highest available and hit 'Next'

Choose deployment type

Deployment types specify common settings for your use case. After creating the domain, you can change these settings at any time.

Deployment type

- Production
Multiple Availability Zones and dedicated master nodes for higher availability.
- Development and testing
One Availability Zone for when you just need an Elasticsearch endpoint.
- Custom
Choose settings from all available options.

Version

Select the version of the Elasticsearch engine for your domain.

Elasticsearch version

6.5

Cancel

Next

Here we provide a name for the domain, and the instance size we would like to use. For our testing purposes, we are going to use a single m5.large. For production purposes, you should always use more than a single node in case there is an issue. You can end up losing data!

Configure domain



A domain is the collection of resources needed to run Elasticsearch. The domain name will be part of your domain endpoint.

Elasticsearch domain name ?

The name must start with a lowercase letter and must be between 3 and 28 characters. Valid characters are a-z (lowercase only), 0-9, and - (hyphen).

Data instances

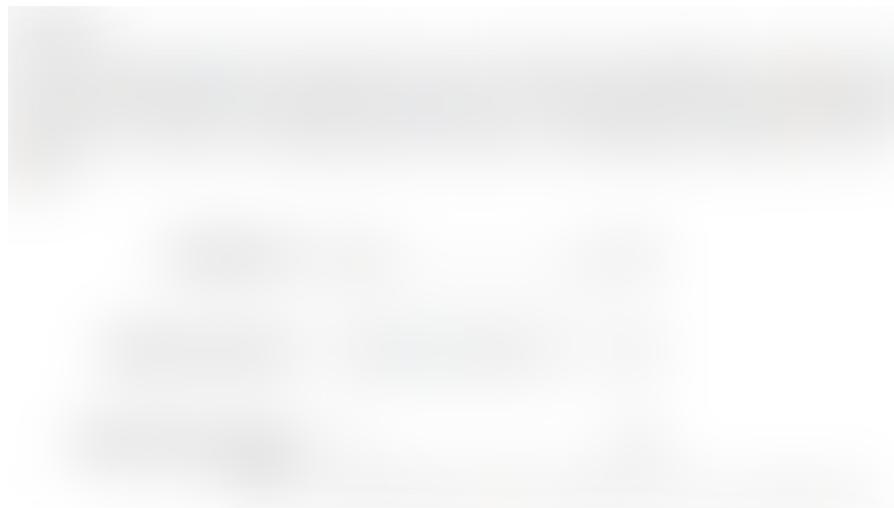
Select an instance type that corresponds to the compute, memory, and storage needs of your application. Consider the size of your Elasticsearch indices, number of shards and replicas, type of queries, and volume of requests. [Learn more](#)

Instance type ?

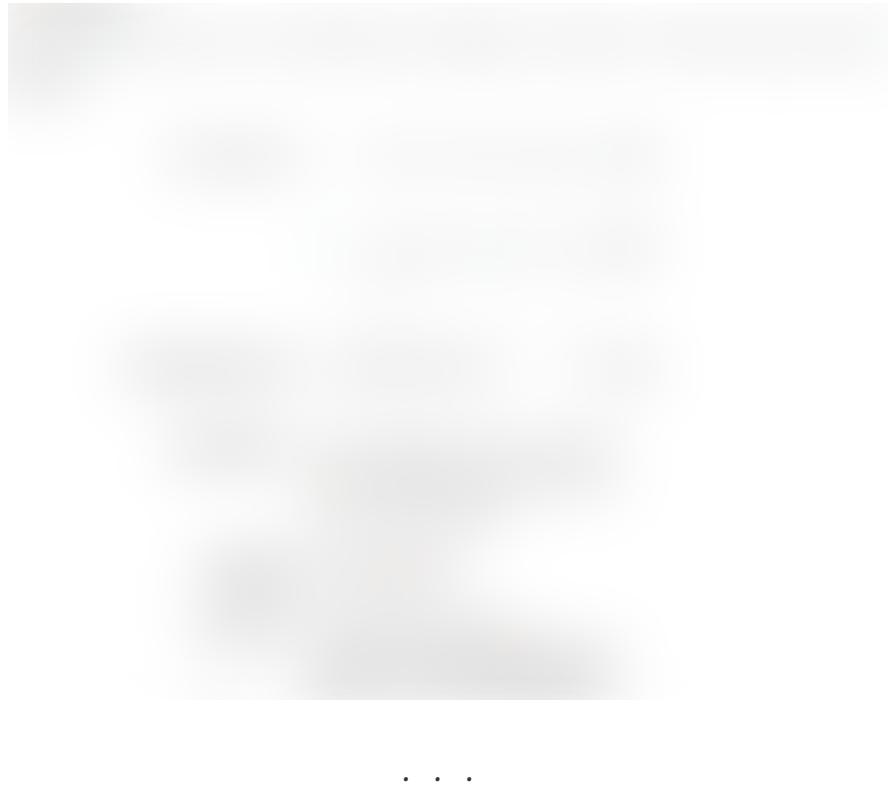
m5.large.elasticsearch instance type needs EBS storage.

Number of instances ?

Next up, set the storage you'd like to provide to ElasticSearch. We will only be loading small amounts of data for testing so you can keep this at 10Gb.



Next up, set the encryption settings for your domain. I'd enable encrypt data at rest just to be on the safe side, I also just use the default key. We only have a single node, so no point enabling node encryption. Hit 'Next' after this.



Next, we have to provide the network settings for our ElasticSearch domain. For testing purposes, we will be using a Public instance restricted to a specific IP address. **YOU SHOULD NOT BE DOING THIS IN PRODUCTION, VPC'S SHOULD BE USED INSTEAD!!!**



Next we have to configure the form of authentication, I don't like using

Cognito so I won't be using it here.



Next we will restrict our domain to just our IP address. To do this, access this link to verify your public IP address —

<http://checkip.amazonaws.com/>

Note it down, and then select from the 'Select a template' dropdown 'Allow access to the domain from specific IP(s)', enter in your IP from the link above and hit OK. You'll see the access policy has been populated with a definition to allow only your IP. This can be modified to a suit a range using the traditional slash notation:

i.e. 10.0.0.0/24 or 10.0.0.0/16 etc.

Mine looks like this (replace XXX.. with your account number). Notice I have also permissioned anything from my account to access the domain, we will need this later for Lambda to access the domain also:

```
1  {
2      "Version": "2012-10-17",
3      "Statement": [
4          {
5              "Effect": "Allow",
6              "Principal": {
7                  "AWS": "arn:aws:iam::XXXXXXXXXXXX:root"
8              },
9              "Action": "es:*",
10             "Resource": "arn:aws:es:eu-west-1:XXXXXXXXXXXX:domain/kinesis/*"
11         },
12         {
13             "Effect": "Allow",
14             "Principal": {
15                 "AWS": "*"
16             },
17             "Action": "es:*",
18             "Resource": "arn:aws:es:eu-west-1:XXXXXXXXXXXX:domain/kinesis/*",
19             "Condition": {
20                 "IpAddress": {
21                     "aws:SourceIp": "1.1.1.1"
22                 }
23             }
24         }
25     ]
26 }
```

```
23      }
24    }
25  ]
26 }
```

[es-domain-access-policy.json](#) hosted with ❤ by GitHub

[view raw](#)

Once that's all done with, hit Next.

• • •

Finally, you'll get a chance to go through all of your chosen settings for the domain. Once you're happy with everything (don't forget things can be changed later) just hit Confirm.

Whilst the domain is being created we'll make a start on the code.

Building the Data Pipeline in Python

First of all, we need to source some data to load. In this guide, I'll be using a sample CSV of Sacramento Crime Records. It's got a few interesting fields you might want to play around with in Kibana, namely location data. You can grab a copy of this data [here](#).

Understanding the different types of data you need to push and how ElasticSearch understands them is very important. If you don't perform the proper data typing, ElasticSearch will only ever understand it as a string. This will become more apparent later when we build the transformation function.



We can see the CSV contains a datetime field we can use in Kibana as our @timefield. We will have to do some work on it first so ElasticSearch can understand it as a datetime.

First of all, I'll always look to build a lightweight script that can quickly collect the required data, perform any transformations required and batch the data up to send to Kinesis. This is where you will discover the largest

bottleneck in time to process data. If you are sourcing data from the results of a SQL query, it's good idea to consider the latency you introduce when connecting to certain data sources and the time it takes for that query to return, as that will affect the time it takes for Kinesis to actually begin processing the data. Here we are using a lightweight example designed for local use, in part 2 we'll look at building a different streaming function so our pipeline is completely serverless with Lambda.

```
1 # necessary imports
2 import boto3
3 import datetime as dt
4 import pandas as pd
5 import time
6
7
8 # function to create a client with aws for a specific service and region
9 def create_client(service, region):
10     return boto3.client(service, region_name=region)
11
12 # function to load data from CSV
13 def load_data(filename):
14     df = pd.read_csv(filename)
15     return df
16
17 # function to correctly display numbers in 2 value format (i.e. 06 instead of 6)
18 def lengthen(value):
19     if len(value) == 1:
20         value = "0" + value
21     return value
22
23 # function for generating new runtime to be used for timefield in ES
24 def get_date():
25
26     today = str(dt.datetime.today()) # get today as a string
27     year = today[:4]
28     month = today[5:7]
29     day = today[8:10]
30
31     hour = today[11:13]
32     minutes = today[14:16]
33     seconds = today[17:19]
34
35     # return a date string in the correct format for ES
36     return "%s/%s/%s %s:%s:%s" % (year, month, day, hour, minutes, seconds)
37
38 # function to modify the date time to be correctly formatted
39 def modify_date(data):
40
41     dates = data['cdatetime'] # get the datetime field
42
43     new_dates = [] # create empty lists
44     load_times = []
45
46     load_time = get_date() # get current time
47
48     # loop over all records
49     for date in dates:
50
51         date = date.replace('/', '-') # replace the slash with dash
```

```

52     date = date + ":00" # add seconds to the datetime
53
54     split = date.split(" ") # split the datetime
55
56     date = split[0] # get just date
57
58     months = date.split('-')[0] # get months
59     days = date.split('-')[1] # days
60     years = "20" + date.split('-')[2] # years
61
62     time = split[1] # get just time
63
64     hours = time.split(':')[0] # get hours
65     minutes = time.split(':')[1] # get minutes
66     seconds = time.split(':')[2] # get seconds
67
68     # build up a string in the right format
69     new_datetime = years + "/" + lengthen(months) + "/" + lengthen(days) + " " + le
70
71     # add it the list
72     new_dates.append(new_datetime)
73     load_times.append(load_time)
74
75     # update the datetime with our transformed version
76     data['cdatetime'] = new_dates
77     data['loadtime'] = load_times
78
79     # return the dataframe
80     return data
81
82     # function for sending data to Kinesis at the absolute maximum throughput
83     def send_kinesis(kinesis_client, kinesis_stream_name, kinesis_shard_count, data):
84
85         kinesisRecords = [] # empty list to store data
86
87         (rows, columns) = data.shape # get rows and columns off provided data
88
89         currentBytes = 0 # counter for bytes
90
91         rowCount = 0 # as we start with the first
92
93         totalRowCount = rows # using our rows variable we got earlier
94
95         sendKinesis = False # flag to update when it's time to send data
96
97         shardCount = 1 # shard counter
98
99         # loop over each of the data rows received
100        for _, row in data.iterrows():
101
102            values = '|'.join(str(value) for value in row) # join the values together by a
103
104            encodedValues = bytes(values, 'utf-8') # encode the string to bytes
105
106            # create a dict object of the row
107            kinesisRecord = {
108                "Data": encodedValues, # data byte-encoded
109                "PartitionKey": str(shardCount) # some key used to tell Kinesis which shard
110            }
111
112
113            kinesisRecords.append(kinesisRecord) # add the object to the list
114            stringBytes = len(values.encode('utf-8')) # get the number of bytes from the st
115            currentBytes = currentBytes + stringBytes # keep a running total

```

```

115     currentBytes = currentBytes + lenBytes # keep a running total
116
117     # check conditional whether ready to send
118     if len(kinesisRecords) == 500: # if we have 500 records packed up, then proceed
119         sendKinesis = True # set the flag
120
121     if currentBytes > 50000: # if the byte size is over 50000, proceed
122         sendKinesis = True # set the flag
123
124     if rowCount == totalRowCount - 1: # if we've reached the last record in the res
125         sendKinesis = True # set the flag
126
127     # if the flag is set
128     if sendKinesis == True:
129
130         # put the records to kinesis
131         response = kinesis_client.put_records(
132             Records=kinesisRecords,
133             StreamName = kinesis_stream_name
134         )
135
136         # resetting values ready for next loop
137         kinesisRecords = [] # empty array
138         sendKinesis = False # reset flag
139         currentBytes = 0 # reset bytecount
140
141         # increment shard count after each put
142         shardCount = shardCount + 1
143
144         # if it's hit the max, reset
145         if shardCount > kinesis_shard_count:
146             shardCount = 1
147
148         # regardless, make sure to incrememnt the counter for rows.
149         rowCount = rowCount + 1
150
151
152     # log out how many records were pushed
153     print('Total Records sent to Kinesis: {}'.format(totalRowCount))
154
155 # main function
156 def main():
157
158     # start timer
159     start = time. time()
160
161     # create a client with kinesis
162     kinesis = create_client('kinesis','eu-west-1')
163
164     # load in data from the csv
165     data = load_data('crime.csv')
166
167     # modify the date and add loadtime field
168     data = modify_date(data)
169
170     # send it to kinesis data stream
171     stream_name = "DataStream"
172     stream_shard_count = 1
173
174     send_kinesis(kinesis, stream_name, stream_shard_count, data) # send it!
175
176     # end timer
177     end = time. time()
178

```

```
179      # log time
180      print("Runtime: " + str(end - start))
181
182  if __name__ == "__main__":
183
184      # run main
185      main()

```

stream-data-kinesis-local.py hosted with ❤ by GitHub [view raw](#)

This script does the following:

- First, we start a timer to capture the script runtime.
- We create a client with Kinesis in the Ireland region (eu-west-1), returns a client
- Loads data in from the CSV renamed to ‘crimes.csv’, returns a pandas DataFrame
- Modifies the DataFrame to correct the datetime field and appends a new datetime field of current time. Returns a DataFrame
- Each record’s fields are joined together using a ‘|’ (pipe) character which we leverage later. We batch up all the data and post it Kinesis to be processed using the client we created, providing the Stream name and shard count. Shards are rotated through to allow for proper fanning out of the workload, however we are only using 1 here.
- Finally, we log the the number of records posted and time taken.

Feel free to modify the variable in main() to suit you.

Running this script loads 7584 records in just under 6 seconds with only one shard.

Obviously, this is a bit of an extreme test and you should be looking to push data more regularly to keep the volume down, this way Kinesis doesn’t get overwhelmed with massive amounts of data all coming in at once resulting in rejected records.

Sourcing the data is by far the most complex of all the steps in this process.

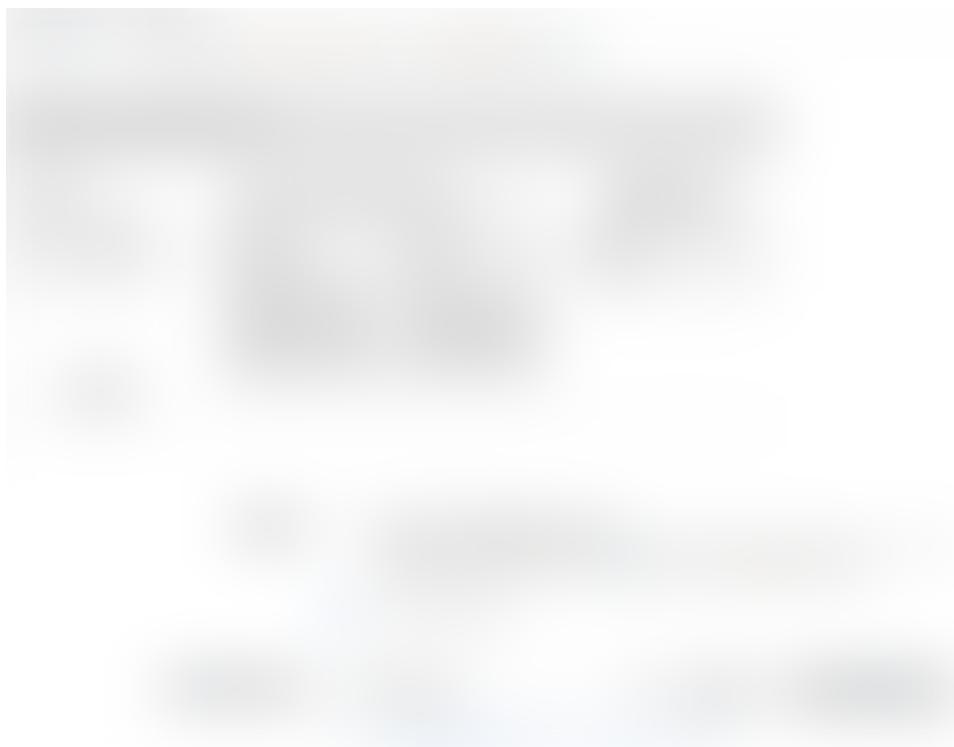
Setting up a Kinesis Firehose with ElasticSearch

In Firehose, hit ‘Create delivery stream’ and provide a name, I will call mine ‘firehose’ for simplicity.



... .

Next, we set the source of our data to be the Kinesis stream we created earlier, then hit Next.



... .

Next we are presented with options for processing records. An important point to note here is that for ElasticSearch to understand the data we have sent up, we must first unpack it into a format it will understand. ElasticSearch understands JSON data so that's what we will need to build.

Currently our data is in the stream waiting to be processed, each record is a pipe delimited string, but we need it as a JSON object. We can use record transformations with a Lambda function to easily do this for us.

Select enabled for record transformation, and let's open a new tab in Lambda.

Building the Transformation Lambda

Here we are going to create a Lambda function that receives data in the form of event JSON containing the records to be transformed. The record's fields are unpacked and then built into a JSON object that ElasticSearch will understand. Here we are using NodeJS 10.x but you may use any language you see fit.

The key to this is understanding the underlying data being passed through the DataStream. It's at this point you will need to perform any data typing ready for the JSON object as you cannot change this once it has been ingested by ES. The keys of your JSON document will also become the fields in ES, so consider the names of them carefully. Another consideration to take into account, is any further processing done here could slow your pipeline down, so I'd recommend doing this before Kinesis receives the data.

You want to get this right as it will affect the uses of the field in Kibana dashboards also.

As everything is a string, we will need to parse certain fields as numbers in order to use them correctly later. Go to Lambda and hit Create Function. Select author from scratch, provide a name (I called mine ‘transform-firehose-data’) and select Node.js 10.x. Let Lambda create a new role as we don’t really need any permissions other than logging. Hit create function.

Once loaded up, we are going to use the following code:

```
1  'use strict';
2
3  console.log('Loading function');
4
5  exports.handler = (event, context, callback) => {
6      let success = 0; // Number of valid entries found
7      let failure = 0; // Number of invalid entries found
8
9      /* Process the list of records and transform them */
10     const output = event.records.map((record) => {
11         // Kinetic data is base64 encoded as decoded base
```

```

11     // Kinesis data is base64 encoded so decode here
12     console.log(record.recordId);
13     const payload = (Buffer.from(record.data, 'base64')).toString('ascii');
14     console.log('Decoded payload:', payload);
15
16     // Split the data into it's fields so we can refer to them by index
17     const match = payload.split('|');
18
19     if (match) {
20         /* Prepare JSON version from Syslog log data */
21         const result = {
22
23             // build all fields from array
24             crime_time: match[0],
25             address: match[1],
26             district: parseInt(match[2]),
27             beat: match[3],
28             grid: parseInt(match[4]),
29             description: match[5],
30             crime_id: parseInt(match[6]),
31             latitude: parseFloat(match[7]),
32             longitude: parseFloat(match[8]),
33             load_time: match[9],
34             location:{
35                 lat: parseFloat(match[7]),
36                 lon: parseFloat(match[8])
37             }
38         };
39         success++;
40         return {
41             recordId: record.recordId,
42             result: 'Ok',
43             data: (Buffer.from(JSON.stringify(result))).toString('base64'),
44         };
45     } else {
46         /* Failed event, notify the error and leave the record intact */
47         failure++;
48         return {
49             recordId: record.recordId,
50             result: 'ProcessingFailed',
51             data: record.data,
52         };
53     }
54 });
55 console.log(`Processing completed. Successful records ${success}, Failed records ${failure}`);
56 callback(null, { records: output });
57 };

```

[transform-data-kinesis.js](#) hosted with ❤ by GitHub [view raw](#)

As you can see we've built the documents fields using the JSON keys, this will be more apparent later on. We perform proper parsing of values into their correct data types. It's important to note that any failures to convert a value will result in the record not arriving in ES. This is more likely due to a field not being completely consistent, you may find some fields that contain a string of 'null' / 'N/A' / 'NaN' to represent a null value when other values are floats or integers.

When this kind of scenario arises, it's good to go back to your original data sourcing/posting script and make the necessary alterations in code so that these issues are corrected. Sometimes this means literally looking through all the rows to fully understand the content of the fields to identify the problematic values. For cases where you do want to have nulls represented correctly in ES, you will have to make changes to the transformation script to check if the field value is equal to your 'null' / 'N/A' / 'NaN' string and then exchange it for an actual null.

Make sure to also increase the timeout of the function to something long like 5 minutes.

Finish creating the Firehose

Hop back to the Firehose creation and select the Lambda function we just created and select the latest version.

We don't need to use record format conversion so we can just leave that disabled. Hit Next when ready.

. . .

Next, we have to select the destination place for our data. We will be using ElasticSearch but you can use anything you wish. You can see the flow this data will take with this setup below:



S3 will take care of anything that fails, great!

Next we are going to configure the ElasticSearch domain and index settings we want to use.

An index can be considered a bit like a schema or a specific table structure. You don't really want to mix different structured data within the same index, generally this leads to a bit of headache when it comes to visualisation as it's not always clear which field are related. I'd suggest as a rule of thumb you use an index per feed/firehose to keep it simple.

A type is yet another way of further classifying data but in the latest version of ElasticSearch posting multiple types to the same indexes is not permitted, so yet again I would go with the same rule as above here.

Next we configure the S3 backup settings and location for any failed records ElasticSearch rejects. This is a good way to monitor that your pipeline is healthy and that records are not failing on their way in. Here we've configured them to go to a specific folder in the bucket.

Hit next when you are ready.

.....

Next, we configure the buffer conditions. This controls how often data is pushed out to the destination. We want this to be as low as possible, set this to 1MB and 60 seconds to decrease the amount of time it takes for data to be output.

.....

Leave everything else default and scroll down to the IAM role section. Here we are going to click ‘Create new or choose’, then allow Firehose to create a new IAM role for this Firehose, by just hitting Allow. Your role should have been created and then you can just click Next.

.....

We get a chance to review all the settings for the Firehose before we create them. When you are ready just hit ‘Create delivery stream’.

Once created, you should see a status of Active shortly after, once you see this your Firehose should be pushing data to ElasticSearch when the DataStream receives data.

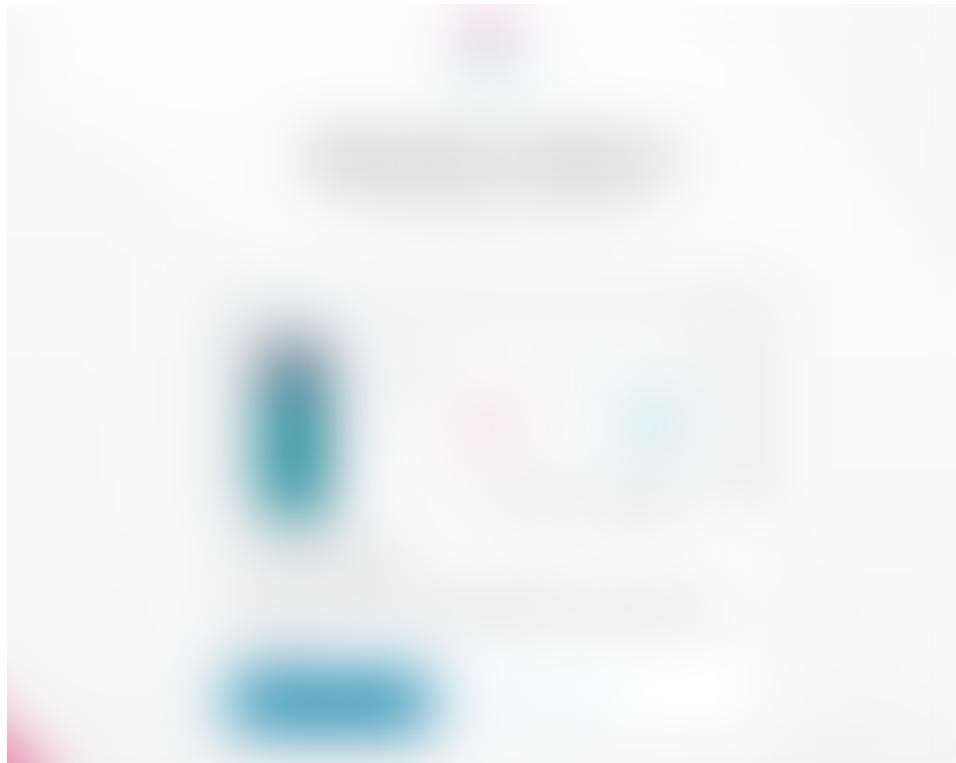
Checking ElasticSearch for Data with Kibana

We can check if ElasticSearch is receiving our data from the Kibana Dashboard.

To access this:

- Go to My Domains, and select your domain.
- Click the Kibana link shown on the Overview tab.

This will launch you into the Kibana UI. If this is your first time accessing Kibana, you'll see something like this:



If you see this, just click 'Explore on my own'.

- Then go to Management > Index Patterns
- Under Define index pattern, enter 'crime-data*'. It should match only one index, and then you can hit next.





On the next screen, we have to select which timefield we want to view the data on.

- If you want to see crimes over time select ‘crime_time’
- If you want every load of data over time select ‘load_time’

I am going for crime_time as I want to see the value used in the data.



Await its creation and then you should see all the fields we captured from the CSV.

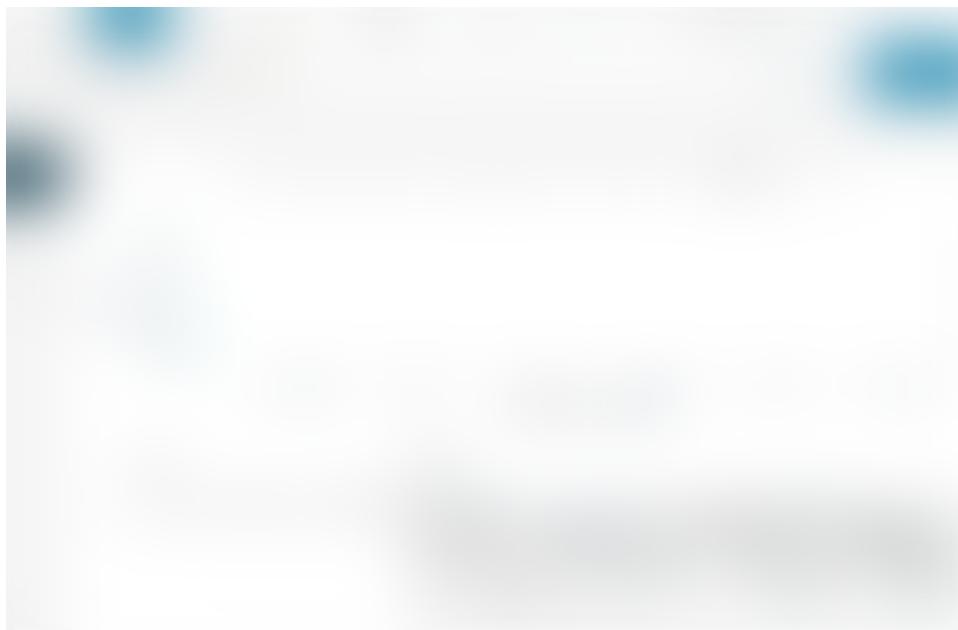
If you want, you can also set this to be the default index for easy access from other screens. If you want to do this hit the star in the top right:





. . .

Going to the Discover tab and setting our time in the top right to view all the way back to 2006 we can see our data has arrived (one batch for each run) in the dashboard. If you used load_time as the time field, you can search for last 1 hour.



You can now leverage all the Kibana tools available to build visualisations and then dashboards etc.

You can learn more about that below:

- Visualisations (start here) — [link](#)
- Dashboards — [link](#)

That's all for Part 1

You should now have data flowing into Kibana whenever you run the collection script. It is just simulation data to show how each of the services interact and how data can end up in the dashboard.

In the next installment, we will be looking at collecting live data from a web service using Lambda as a serverless approach to data streaming and visualisation.

Follow me to know when it's ready! 😊

Discover Medium

Welcome to a place where words matter. On Medium, smart voices and original ideas take center stage - with no ads in sight. Watch

Make Medium yours

Follow all the topics you care about, and we'll deliver the best stories for you to your homepage and inbox. Explore

Explore your membership

Thank you for being a member of Medium. You get unlimited access to insightful stories from amazing thinkers and storytellers. Browse