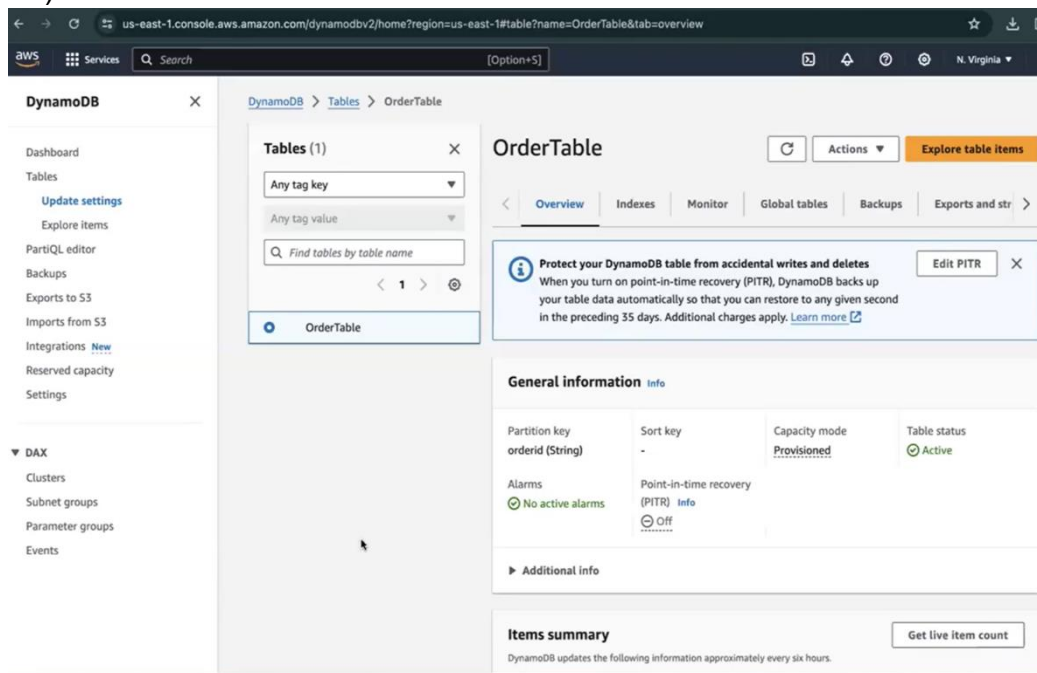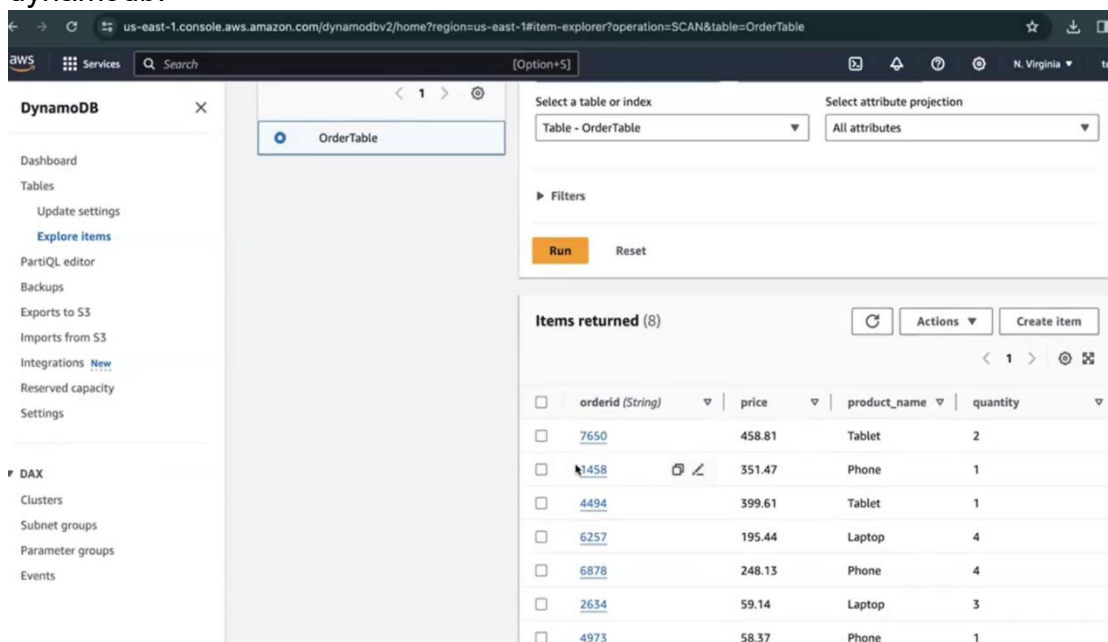# Sales Data Projection

1. this was a sales data streaming project for processing the data from Dynamo db to destination S3 target path.
2. there was some source from where records(orderid, product_name, quantity, price) were being written into the nosql dynamodb(which is a key value based db).



3. we created a dynamodb table where we targeted particular partition key(orderid) and the entire data(orderid, product_name, quantity, price) used to get ingested as a value there. That's how we were publishing a record in dynamodb.

4. next our task was to capture changes happening in table. Records were getting inserted, updated/deleted to table in real time. CDC is whatever changes are happening in real time we capture those changes. That's what CDC is. So we needed to enable dynamodb streams from 'Exports & Streams' option to perform CDC. As we were able to consume the CDC changes, hence created CDC pipeline.

5. then we set up a kinesis stream which is a queuing mechanism to hold our data. Kinesis stream has a concept of shards just like publishing data to partitions. Data here arrives at kinesis stream in different shard.
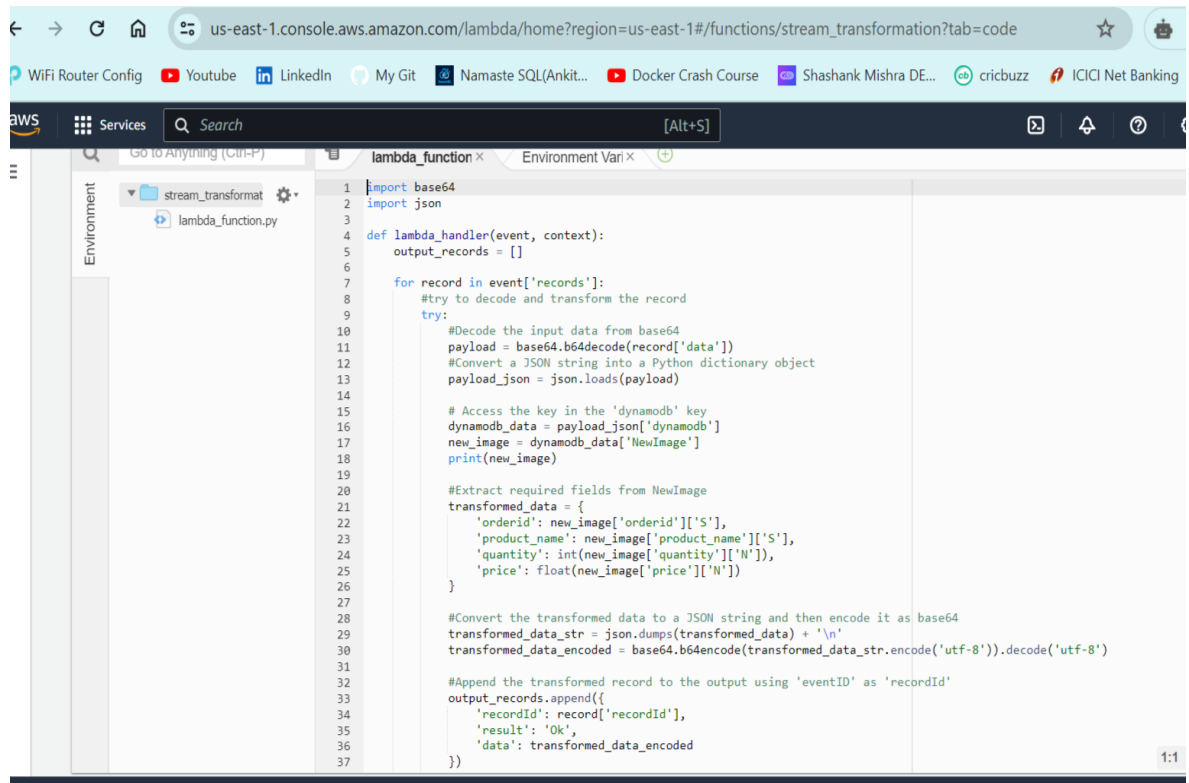
6. then we had a event bridge pipe in between dynamodb stream and kinesis stream. We set a batch size consisting of no of message to be sent to any kinesis stream shard. Set up required IAM role to eventbridge pipe to access dynamo stream & kinesis stream. Using event bridge we used to flow data to kinesis stream.



7. next we created a firehose stream which is a delivery stream. It processed data(like ingest, transform) and delivered streaming data to destinations(like data lakes, data warehouses). Here source is kinesis stream and destination is s3.
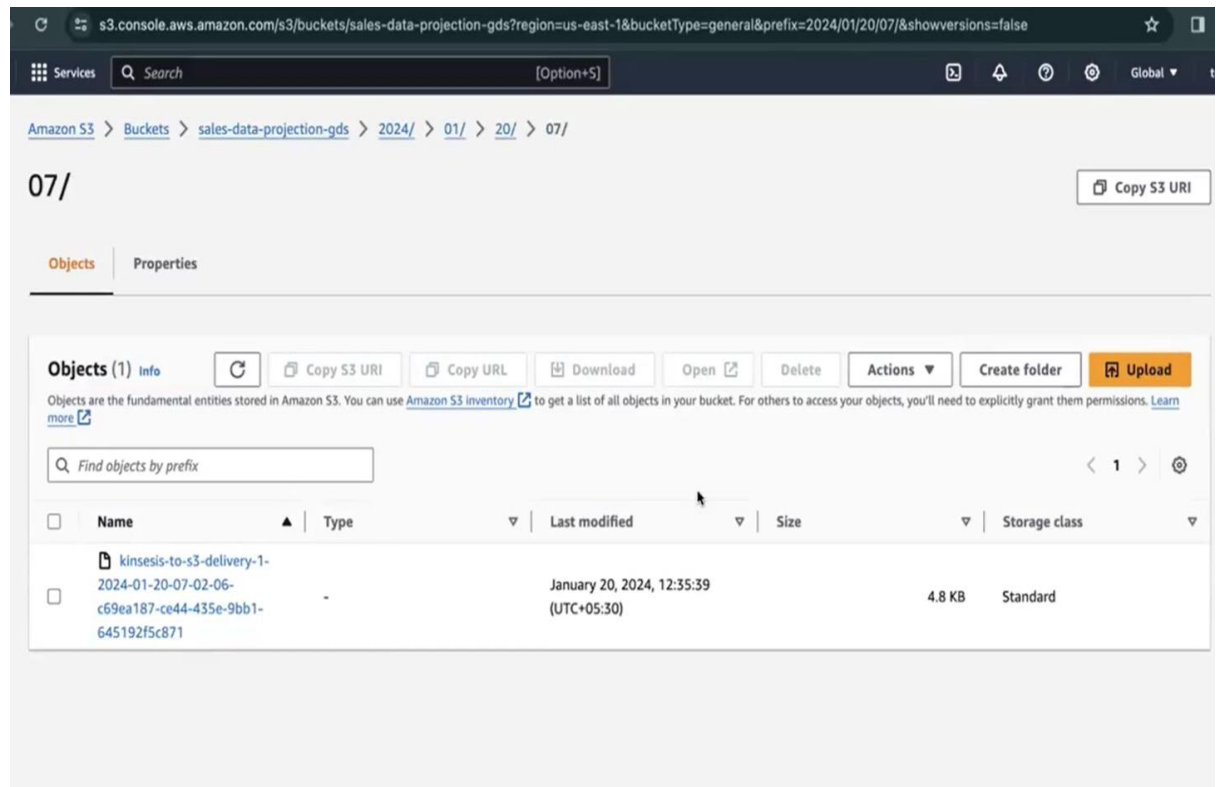
8. using delivery stream we were batching the records and transforming those mini batch records with the help of lambda function.



```python
import base64
import json

def lambda_handler(event, context):
    output_records = []

    for record in event['records']:
        #try to decode and transform the record
        try:
            #Decode the input data from base64
            payload = base64.b64decode(record['data'])
            #Convert a JSON string into a Python dictionary object
            payload_json = json.loads(payload)

            # Access the key in the 'dynamodb' key
            dynamodb_data = payload_json['dynamodb']
            new_image = dynamodb_data['NewImage']
            print(new_image)

            #Extract required fields from NewImage
            transformed_data = {
                'orderid': new_image['orderid']['S'],
                'product_name': new_image['product_name']['S'],
                'quantity': int(new_image['quantity']['N']),
                'price': float(new_image['price']['N'])
            }

            #Convert the transformed data to a JSON string and then encode it as base64
            transformed_data_str = json.dumps(transformed_data) + '\n'
            transformed_data_encoded = base64.b64encode(transformed_data_str.encode('utf-8')).decode('utf-8')

            #Append the transformed record to the output using 'eventID' as 'recordId'
            output_records.append({
                'recordId': record['recordId'],
                'result': 'Ok',
                'data': transformed_data_encoded
            })
```

9. Once this is done, then delivering the transformed batch records to S3.

10. Lastly, we created crawler over s3 target path which resulting in catalog table creation. Also, needed to add classifier for json format(i.e. $.columns).



11. using athena query we then did analysis