

Change Data Capture (CDC)

What is Change Data Capture? Example and tools to do Change Data Capture

Change Data Capture (CDC) is a technique used to capture changes made to a database so that they can be propagated to other systems in near real-time. CDC is often used in data integration, data replication, and data warehousing scenarios to ensure that the target systems have up-to-date and accurate data.

CDC works by monitoring the transaction log of a database and capturing changes made to the data. The captured changes can then be transformed and loaded into another system, such as a data warehouse or reporting database. CDC can capture changes at the row level, which means that only the changed data is propagated, reducing the amount of data transferred and increasing efficiency.

There are several tools and technologies that can be used to implement CDC, including:

1. Microsoft SQL Server Change Data Capture - This feature is built into Microsoft SQL Server and can be used to capture and propagate changes made to SQL Server databases. It provides a reliable and efficient way to capture changes and propagate them to other systems.
2. Oracle GoldenGate - This tool is used to capture and replicate changes made to Oracle databases. It provides real-time data integration and can capture changes made to databases on-premises or in the cloud.
3. Apache Kafka - Kafka is a distributed streaming platform that can be used for real-time data integration, including CDC. It can capture changes made to databases and propagate them to other systems in real-time, making it a popular choice for data streaming use cases.
4. Attunity Replicate - This tool is used for data integration, replication, and data warehousing scenarios. It can capture changes made to databases and propagate them to other systems in near real-time, and it supports a wide variety of data sources and targets.

An example of using CDC might be capturing changes made to a customer database, such as new orders or changes to customer details, and propagating those changes to a data warehouse for analysis and reporting purposes. CDC can also be used to keep databases in sync across different systems or to replicate data across different environments.

What programming language to use in CDC?

The programming language used to implement CDC (Change Data Capture) depends on the technology being used to capture and process changes made to a database.

For example, if you are using the built-in CDC feature in Microsoft SQL Server, you can use SQL scripts or T-SQL (Transact-SQL) code to capture changes and process them.

If you are using a tool like Apache Kafka, you can use programming languages such as Java, Python, or Scala to implement the CDC solution.

Similarly, if you are using a cloud-based CDC solution like AWS DMS or Google Cloud Dataflow, you can use programming languages such as Java, Python, or Go to implement the solution.

Overall, the choice of programming language depends on the technology being used and the specific requirements of the CDC

Techniques to implement CDC?

There are several techniques to implement Change Data Capture (CDC) in a data integration process. Here are a few commonly used techniques:

1. **Database Triggers:** One of the simplest ways to implement CDC is through the use of database triggers. A trigger can be set up on a table to automatically capture any changes to the data in the table. When a change occurs, the trigger can be configured to write the change details to a separate CDC table or a change file.
2. **Change Tracking:** Many relational database management systems (RDBMS) provide built-in mechanisms to track changes to data. For example, Microsoft SQL Server has a feature called Change Tracking that captures changes to selected tables in the database. Change Tracking can be used to identify new, updated, and deleted rows.
3. **Log-based CDC:** Log-based CDC involves capturing changes to a database through the transaction log. This method requires a more complex setup but offers more robust and reliable CDC. Log-based CDC can be implemented using tools such as Apache Kafka, Debezium, or Oracle GoldenGate.
4. **Third-party Tools:** There are many third-party tools available in the market that can be used to implement CDC. These tools offer features such as automated change capture, data mapping, data transformation, and data delivery. Some examples of CDC tools are Attunity Replicate, HVR, and Qlik Replicate.

The choice of technique for implementing CDC will depend on factors such as the type of data source, the complexity of the integration process, the level of data accuracy and completeness required, and the budget and resources available.

CDC setup in AWS

Change Data Capture (CDC) can be used in AWS S3 and Google BigQuery to capture changes made to databases and propagate them to other systems.

In both AWS S3 and Google BigQuery, CDC can be used to capture changes made to databases in real-time and propagate those changes to other systems. This can be useful for data warehousing, analytics, and reporting scenarios where up-to-date data is required for analysis and decision-making.

AWS S3 - To use CDC in AWS S3, you can use AWS Database Migration Service (DMS). DMS can capture changes made to databases, including AWS RDS, Aurora, and other databases, and replicate those changes to S3. To set up CDC with DMS, you will need to create a replication instance, source endpoint, and target endpoint. You can then configure DMS to capture changes made to the database using the CDC feature. The changes can be stored in S3 in Apache Parquet format, which is a columnar storage format optimized for analytics workloads.

👍 If you use Google BigQuery - To use CDC in Google BigQuery, you can use a tool like Alooma or Stitch Data. These tools can capture changes made to databases, including PostgreSQL, MySQL, and MongoDB, and replicate those changes to BigQuery. To set up CDC with Alooma or Stitch Data, you will need to create a data source connection to your database, select the tables to replicate, and configure the CDC settings. The changes can be stored in BigQuery in a format optimized for analytics workloads, such as columnar or partitioned formats.

Example of CDC in Python in AWS S3 using DMS (MySQL)

Here's an example of how to set up CDC in Python using AWS DMS and S3:

1. First, you'll need to set up AWS DMS and configure the source and target endpoints. In this example, we'll assume you have an AWS DMS endpoint set up to capture changes from a MySQL database.
2. Next, you'll need to create an S3 bucket to store the CDC data. You can use the AWS Management Console or the AWS CLI to create the bucket. Or can implement this example code

```

1 import boto3
2
3 # set up boto3 client for S3
4 s3 = boto3.client('s3', region_name='us-east-1')
5
6 # create an S3 bucket for storing CDC data
7 bucket_name = 'my-cdc-bucket'
8 s3.create_bucket(Bucket=bucket_name)
9

```

3. Once your bucket is set up, we need to set up a DMS replication instance and endpoint, and then create a CDC task. You can configure AWS DMS to write the CDC data to S3. In the AWS DMS console, go to the task settings and select the "Target" tab. Here, you can select "Amazon S3" as the target endpoint and configure the S3 settings (e.g. bucket name, folder structure, etc.).
4. After you've configured the S3 settings, you can start the data migration task to begin capturing changes from your MySQL database and writing them to S3.
5. Once the changes are captured, you can use Python to query and process the data in S3. Here's some sample code to get started:

```

1 import boto3
2 import pandas as pd
3 # Initialize a client object
4 s3 = boto3.client('s3', region_name='us-east-1')
5
6 # Define the bucket and folder names
7 bucket_name = 'my-bucket'
8 folder_name = 'my-folder'
9
10 # List all the objects in the folder
11 objects = s3.list_objects_v2(Bucket=bucket_name, Prefix=folder_name)['Contents']
12
13 # Loop through each object and load the CDC data into a dataframe
14 dfs = []
15 for obj in objects:
16     obj_key = obj['Key']
17     obj_body = s3.get_object(Bucket=bucket_name, Key=obj_key)['Body']
18     df = pd.read_csv(obj_body, delimiter=',')
19     dfs.append(df)
20
21 # Concatenate all the dataframes into one
22 df = pd.concat(dfs)
23

```

```

24 # Filter the data by date range
25 start_date = '2022-01-01'
26 end_date = '2022-03-01'
27 df = df[(df['timestamp'] >= start_date) & (df['timestamp'] < end_date)]
28

```

This code creates a client object to connect to S3, defines the bucket and folder names, and queries the folder to load the CDC data into dataframes. It then concatenates all the dataframes into one and filters the data by a date range.

You can then use Python libraries like Pandas and Matplotlib to analyze and visualize the data as needed.

Note that this is just a simple example and the specific implementation will depend on your AWS DMS settings, S3 bucket structure, and other factors.

Other way, Here is the explanation code that combine step 3,4,5 together : We create a replication instance, subnet group, and endpoint using the

`create_replication_instance`, `create_replication_subnet_group`, and `create_endpoint` methods, respectively. We then create a CDC task using the `create_replication_task` method, specifying the source and target endpoints, replication instance, and table mappings file. The `CdcStartPosition` parameter is set to `TRUNCATE_BEFORE_UPDATE` to truncate the target table before updates, and `CdcStopPosition` is set to `LAST_RECORD` to stop the task after processing the last record in the source database.

```

1 import boto3
2
3 # set up boto3 client for DMS
4 dms = boto3.client('dms')
5
6 # create a replication instance
7 replication_instance_id = 'my-replication-instance'
8 dms.create_replication_instance(
9     ReplicationInstanceIdentifier=replication_instance_id,
10     AllocatedStorage=100,
11     EngineVersion='3.3.1',
12     ReplicationInstanceClass='dms.t2.micro',
13     VpcSecurityGroupIds=['sg-12345678'],
14     AvailabilityZone='us-west-2a',
15     PubliclyAccessible=False

```

```
16 )
17
18 # create a replication subnet group
19 subnet_group_id = 'my-subnet-group'
20 dms.create_replication_subnet_group(
21     ReplicationSubnetGroupIdentifier=subnet_group_id,
22     ReplicationSubnetGroupDescription='My replication subnet group',
23     SubnetIds=['subnet-12345678', 'subnet-23456789']
24 )
25
26 # create a replication endpoint
27 endpoint_id = 'my-endpoint'
28 dms.create_endpoint(
29     EndpointIdentifier=endpoint_id,
30     EndpointType='source',
31     EngineName='aurora',
32     EngineDisplayName='Amazon Aurora',
33     Username='my-username',
34     Password='my-password',
35     ServerName='my-aurora-cluster.cluster-abcdefg.us-west-2.rds.amazonaws.com',
36     Port=3306,
37     DatabaseName='my-database',
38     SslMode='none',
39     ReplicationInstanceIdentifier=replication_instance_id,
40     Tags=[
41         {'Key': 'environment', 'Value': 'prod'}
42     ]
43 )
44
45 # create a CDC task
46 task_id = 'my-cdc-task'
47 dms.create_replication_task(
48     ReplicationTaskIdentifier=task_id,
49     SourceEndpointArn='arn:aws:dms:us-west-2:123456789012:endpoint:my-endpoint',
50     TargetEndpointArn='arn:aws:dms:us-west-2:123456789012:endpoint:s3-target',
51     ReplicationInstanceArn='arn:aws:dms:us-west-2:123456789012:rep:my-replicatio
52     MigrationType='cdc',
53     TableMappings='file://table-mappings.json',
54     CdcStartPosition='TRUNCATE_BEFORE_UPDATE',
55     CdcStopPosition='LAST_RECORD',
56     Tags=[
57         {'Key': 'environment', 'Value': 'prod'}
58     ]
59 )
60
```


Code implement for CDC in Python in AWS using AWS Lambda and CloudWatch (MySQL)

1. Create an Amazon RDS instance and a MySQL database.
2. Install the necessary Python libraries: boto3, mysql-connector-python, and logzero.

```
1 pip install boto3 mysql-connector-python logzero
```

3. Configure AWS credentials using boto3. You can use the following code to set up your credentials

```
1 import boto3
2 session = boto3.Session(profile_name='myprofile')
3 client = session.client('rds')
4
```

4. Create an S3 bucket to store the CDC logs

```
1 import boto3
2 session = boto3.Session(profile_name='myprofile')
3 s3 = session.client('s3')
4
5 bucket_name = 'my-cdc-bucket'
6 s3.create_bucket(Bucket=bucket_name)
```

5. Create a Lambda function that reads the changes made to the database and writes them to an S3 bucket.

```
1 import json
2 import boto3
3 import mysql.connector
```

```

4  from logzero import logger
5
6  def lambda_handler(event, context):
7      # Get database connection details from environment variables
8      host = os.environ['DB_HOST']
9      user = os.environ['DB_USER']
10     password = os.environ['DB_PASSWORD']
11     database = os.environ['DB_NAME']
12
13     # Create connection to the database
14     connection = mysql.connector.connect(
15         host=host,
16         user=user,
17         password=password,
18         database=database
19     )
20
21     # Get a cursor to execute queries
22     cursor = connection.cursor()
23
24     # Get the latest change timestamp from the S3 bucket
25     s3 = boto3.client('s3')
26     latest_timestamp = s3.get_object(
27         Bucket=os.environ['CDC_BUCKET'],
28         Key='latest_timestamp.txt' )['Body'].read().decode('utf-8')
29
30     # Get the changes made to the database since the latest timestamp
31     cursor.execute(f"SELECT * FROM mytable WHERE modified_date > '{latest_times
32 changes = cursor.fetchall()
33
34     # Write the changes to an S3 bucket
35     s3.put_object(
36         Bucket=os.environ['CDC_BUCKET'],
37         Key=f'{datetime.now().strftime("%Y-%m-%d-%H-%M-%S")}.json',
38         Body=json.dumps(changes)
39     )
40
41     # Update the latest timestamp in the S3 bucket
42     s3.put_object(
43         Bucket=os.environ['CDC_BUCKET'],
44         Key='latest_timestamp.txt',
45         Body=datetime.now().strftime("%Y-%m-%d %H:%M:%S")
46     )
47
48     logger.info(f"{len(changes)} changes made to the database")

```


6. Create a CloudWatch Events rule to trigger the Lambda function every time a change is made to the database.

```
1 import boto3
2 session = boto3.Session(profile_name='myprofile')
3 events = session.client('events')
4
5 rule_name = 'my-cdc-rule'
6 function_arn = 'arn:aws:lambda:us-east-1:123456789012:function:my-cdc-function'
7
8 events.put_rule(
9     Name=rule_name,
10    ScheduleExpression='cron(0/5 * * * ? *)', # run every 5 minutes
11    State='ENABLED'
12 )
13
14 events.put_targets(
15     Rule=rule_name,
16     Targets=[
17         {
18             'Id': 'my-cdc-target',
19             'Arn': function_arn
20         }
21     ]
22 )
23
```

This code sets up a CloudWatch Events rule that triggers the Lambda function every 5 minutes. The Lambda function then reads the changes made to the database from the last run and writes them to an S3 bucket.

CDC in AWS using the boto3 library to connect to an Amazon RDS PostgreSQL database.

First, we need to set up an Amazon RDS PostgreSQL database and enable CDC for the desired tables. This can be done through the AWS Management Console or using the AWS CLI. Once CDC is enabled, we can use the boto3 library to create a connection to the database and retrieve the changes.

Here is an example code. This code uses the boto3 library to retrieve the database endpoint, port, and credentials from AWS, and the psycopg2 library to connect to the database and retrieve the changes from the CDC enabled tables. The changes can then be processed as needed.

```
1 import boto3
2 import psycopg2
3 from psycopg2 import sql
4
5 # Set up connection to AWS RDS PostgreSQL database
6 client = boto3.client('rds')
7 response = client.describe_db_instances(DBInstanceIdentifier='my-db-instance')
8 db_endpoint = response['DBInstances'][0]['Endpoint']['Address']
9 db_port = response['DBInstances'][0]['Endpoint']['Port']
10 db_user = response['DBInstances'][0]['MasterUsername']
11 db_password = 'my-db-password'
12
13 # Connect to the database
14 conn = psycopg2.connect(host=db_endpoint, port=db_port, dbname='my-db', user=db
15 cursor = conn.cursor()
16
17 # Get list of CDC enabled tables
18 query = sql.SQL("SELECT table_name FROM information_schema.tables WHERE table_c
19 cursor.execute(query)
20 tables = cursor.fetchall()
21
22 # Retrieve changes from CDC enabled tables
23 for table in tables:
24     query = sql.SQL("SELECT * FROM {} WHERE xmin::text::bigint > %s").format(sq
25     cursor.execute(query, [0])
26     changes = cursor.fetchall()
27
28     # Process changes
29     for change in changes:
30         # Process change
31         pass
32
33 cursor.close()
34 conn.close()
35
```

Bonus: Example of CDC in Python in Google Big Query using Stitch Data

Here's an example of how to set up CDC in Python using Stitch Data and Google BigQuery:

1. First, you'll need to sign up for Stitch Data and create a new data pipeline.
2. Next, you'll need to select your data source (e.g. MySQL, PostgreSQL, MongoDB, etc.) and configure the connection settings.
3. Once your data source is connected, you can configure the CDC settings. In the Stitch Data dashboard, go to the "Advanced" tab and select "CDC" to configure the CDC settings. Here, you can specify the tables to capture changes for, the frequency of updates, and other settings.
4. After you've configured the CDC settings, you can start the data pipeline to begin capturing changes from your database.
5. Once the changes are captured, you can use Python to query and process the data in Google BigQuery. Here's some sample code to get started:

```
1 from google.cloud import bigquery
2 # Initialize a client object
3 client = bigquery.Client()
4
5 # Define the dataset and table names
6 dataset_id = 'my_dataset'
7 table_id = 'my_table'
8 # Create a reference to the table
9 table_ref = client.dataset(dataset_id).table(table_id)
10 # Define the query to load the CDC data into BigQuery
11 query = f"""
12     SELECT * FROM `{table_ref}`
13     WHERE _sdc_batched_at >= TIMESTAMP_SUB(CURRENT_TIMESTAMP(), INTERVAL 1 DAY)
14     """
15
16 # Execute the query and save the results to a dataframe
17 df = client.query(query).to_dataframe()
18
```

This code creates a client object to connect to Google BigQuery, defines the dataset and table names, and queries the table to load the CDC data into a dataframe.

You can then use Python libraries like Pandas and Matplotlib to analyze and visualize the data as needed.

Note that this is just a simple example and the specific implementation will depend on your data source, CDC settings, and other factors.

--

What alternatives of CDC in AWS, GCP and Azure?

There are several alternatives to Change Data Capture (CDC) available in cloud platforms like AWS, GCP, and Azure. Here are some of the alternatives:

1. **AWS DMS (Database Migration Service):** AWS DMS is a managed service that can be used to migrate data from one database to another. It supports both homogeneous and heterogeneous migrations, and also supports ongoing replication for CDC.
2. **AWS Kinesis Data Streams:** AWS Kinesis Data Streams is a real-time data streaming service that allows you to capture and process data in real time. You can use Kinesis Data Streams to ingest data from various sources and stream it to AWS services like S3, Redshift, or Elasticsearch.
3. **GCP Dataflow:** GCP Dataflow is a fully managed service for building batch and streaming data pipelines. It provides a unified programming model for both batch and streaming data processing and supports CDC.
4. **GCP Cloud Functions:** GCP Cloud Functions is a serverless compute service that allows you to run your code in response to events. You can use Cloud Functions to create a serverless CDC pipeline that triggers changes to your source database and writes the changes to your target database.
5. **Azure Stream Analytics:** Azure Stream Analytics is a fully managed service for real-time data stream processing. You can use Stream Analytics to ingest data from various sources, process the data using SQL-like queries, and output the results to Azure services like Blob storage or SQL Database.
6. **Azure Functions:** Azure Functions is a serverless compute service that allows you to run your code in response to events. You can use Functions to create a serverless CDC pipeline that triggers changes to your source database and writes the changes to your target database.

These are just a few examples of the CDC alternatives available in cloud platforms. The choice of tool depends on your specific use case and requirements.