# AI-Powered Video Surveillance System with Chatbot-Based Querying and Machine Learning Pipelines on AWS

## 1. Data Ingestion

The data ingestion layer is responsible for efficiently capturing live video streams from multiple sources, such as security cameras deployed across various locations, and routing these streams into the AWS cloud for analysis, storage, and querying. For this project, AWS Kinesis Video Streams (KVS) was selected due to its ability to handle real-time video streams at scale, offering seamless integration with other AWS services like Rekognition, Lambda, and S3.

In addition to handling the video streams, we will also attach important metadata as KVS embedded tags to each stream to provide critical context such as the location of the camera, the camera ID, and other relevant details.

1. **Camera ID (`camera_id`)**
2. **Camera Location (`location`)**
3. **Video Resolution (`resolution`)**
4. **Frame Rate (`frame_rate`)**
5. **Timestamp (`timestamp`)**
6. **Stream ARN (`stream_arn`)**
   - The Amazon Resource Name (ARN) of the Kinesis Video Stream that is ingesting the video. This is a unique identifier for the stream in the AWS ecosystem.
   - Example: `"arn:aws:kinesisvideo:region:account-id:stream/stream-name"`
7. **Video Format (`video_format`)**:
   - The format in which the video stream is captured, such as H.264, MJPEG, or H.265.

## 2. Storage

Amazon S3 serves as the storage solution for video footage from each client's security cameras. The video files will be uploaded to S3 via Lambda, with client-specific folder structures in place to ensure proper organization and access control. The S3 bucket will be configured to trigger an AWS Lambda function when a new video file is uploaded. This Lambda function is responsible for initiating the preprocessing and sending data to the Rekognition API for video analysis.

`datetime`: This module is used to timestamp video files when they are saved in Amazon S3.

`boto3`: This is the AWS SDK for Python. It allows Python developers to interact with AWS services like Kinesis, S3, Rekognition, Lambda, etc. Here, it's being used to interact with Amazon Kinesis Video Streams and Amazon S3. With Boto3, you can perform tasks like creating an S3 bucket, uploading files, launching EC2 instances, retrieving data from DynamoDB, and more, all through Python code. The Boto3 Client API allows you to interact directly with AWS services using low-level API calls to create streams, ingest media, and retrieve data from Kinesis streams.

`kinesis_client = boto3.client('kinesisvideo')`: This line creates a client for the Amazon Kinesis Video Streams service, which will allow the code to interact with Kinesis for retrieving video streams.

`s3_client = boto3.client('s3')`: This client for Amazon S3 is used to upload the video data to a specified S3 bucket.

**Setting Up Event-Driven Triggers:**
S3 event notifications are used to automatically trigger Lambda functions whenever a new file is uploaded to the S3 bucket. This is the first step in the event-driven pipeline.

- In the AWS S3 console, configure the S3 bucket to generate an event notification on ObjectCreated events (whenever a new video is uploaded). The event will trigger an AWS Lambda function that processes the video.

## 3. Define the Lambda Handler Function

```
def lambda_handler(event, context):
```

This is the Lambda handler function, the entry point for an AWS Lambda function. Lambda automatically invokes this function when triggered by a Kinesis video stream event.

- `event`: Contains the data passed by the Kinesis video stream being processed, including the `streamARN`. The `streamARN` in the event object tells the Lambda function which specific video stream triggered the function. Since we have multiple Kinesis video streams in our AWS environment, the `streamARN` is used to identify which specific  stream the Lambda function needs to process and retrieve media from. ARNs are used across AWS to uniquely identify resources such as S3 buckets, EC2 instances, Lambda functions, and Kinesis streams. It contains the ARN prefix, the AWS resource namespace, the service namespace (in this case, Kinesis Video Streams), the AWS region where the stream is located, the AWS account ID of the stream owner, the name of the video stream, and the stream timestamp.
- `context`: Provides runtime information about the Lambda function's execution environment. It contains metadata about the Lambda invocation, including the name of the Lambda function, the amount of memory allocated to the Lambda function, a unique identifier for the Lambda invocation, and information on where logs are stored in Amazon CloudWatch.

## Why Use AWS Lambda?

AWS Lambda is chosen in this case due to its serverless nature, automatic scaling, cost efficiency, and event-driven architecture. Here's why Lambda is ideal for this project as opposed to other AWS compute options:

- **No Server Management**: Lambda is fully serverless, which means you don't have to manage any infrastructure. AWS automatically handles the infrastructure, allowing you to focus solely on the business logic (processing video streams, storing them, etc.).
- **Event-Driven**: Lambda automatically triggers when an event occurs (like a new video stream or an upload to S3). This makes Lambda perfect for event-driven architectures, such as the one required in this project where the function is triggered when a new video stream is available.
- **Auto-Scaling**: Lambda automatically scales the compute resources up or down based on the number of incoming events (like video streams). If multiple Kinesis video streams are incoming simultaneously, Lambda will automatically scale the number of function invocations to handle the load.

- ○ **Concurrent Execution**: Lambda can handle hundreds of thousands of concurrent executions, meaning it's well-suited for processing multiple video streams or requests at the same time, without needing to manually configure scaling parameters.
- ○ **Pay-as-You-Go**: With Lambda, you only pay for the compute time you use, which makes it extremely cost-effective for event-driven tasks like this one. You are billed based on the number of function invocations and the time it takes for the function to run, rounded to the nearest millisecond.
- ○ **No Idle Costs**: Since Lambda is serverless, you don't pay for idle time. If no video streams are being processed, you're not incurring costs, which is highly efficient compared to constantly running servers or containers.
- ○ **Native AWS Integration**: Lambda integrates seamlessly with other AWS services, such as Kinesis Video Streams, S3, Rekognition, and DynamoDB. In this case, Lambda can easily respond to new video streams from Kinesis, process them, and store them in S3 without needing custom setup or configuration for service-to-service communication.
- ○ **Event Triggers**: Lambda can be triggered by various AWS services, including S3 (file uploads), DynamoDB (changes in data), and Kinesis (new video fragments). In your use case, it's triggered when a video stream event happens, providing an instant response to real-time video streams.
- ○ **No Maintenance Required**: Since Lambda is fully managed, there is no need for server maintenance, operating system patching, or infrastructure monitoring. AWS automatically handles all the infrastructure components, so you only need to focus on the Lambda function's code.
- ○ **Automatic Updates**: AWS continuously updates and manages Lambda infrastructure behind the scenes, including security patches, without any downtime for your applications.
- ○ **Low-Latency Execution**: Lambda offers low-latency execution times, especially for real-time video ingestion tasks like the one described in your project. Since it's event-driven, the function is invoked instantly when the event occurs, processing video streams in near-real-time.
- ○ **Cold Starts**: While Lambda functions can experience cold starts (a slight delay when the function is invoked for the first time), for your use case (processing video streams periodically), this delay will be minimal and not a significant issue.

## 4. Extract the Kinesis Video Stream ARN

- ● `stream_arn = event['streamARN']`: This extracts the Amazon Resource Name (ARN) of the video stream from the incoming `event`. This value is passed in the event, which the Lambda function receives when triggered by a Kinesis event.

## 5. Retrieve the Kinesis Data Endpoint

```
data_endpoint = kinesis_client.get_data_endpoint(
    StreamARN=stream_arn,
    APIName='GET_MEDIA')['DataEndpoint']
```

Before fetching video data from the Kinesis stream, the client needs to get a data endpoint. The data endpoint is the URL from which the actual media (video stream) will be retrieved.

- ● **Function**: `get_data_endpoint` is a method of the `kinesisvideo` client. It requests the appropriate data endpoint for retrieving media from the video stream.

- **Parameters**:
  - `StreamARN=stream_arn`: This specifies which Kinesis video stream to retrieve the endpoint from.
  - `APIName='GET_MEDIA'`: This specifies that you want an endpoint for retrieving media data. GET_MEDIA is the API operation for fetching the video stream from Kinesis.

## 6. Create a Media Client for Kinesis Video Streams

```
media_client = boto3.client('kinesis-video-media', endpoint_url=data_endpoint)
```

This line creates a media client to interact with Kinesis Video Streams using the media-specific data endpoint. Since the actual video media is stored at a different endpoint (returned by `get_data_endpoint`), the `kinesis-video-media` client is used to communicate with this endpoint.

## 7. Fetch Video Data from the Kinesis Stream

```
video_data = media_client.get_media(
    StreamARN=stream_arn,
    StartSelector={'StartSelectorType': 'NOW'})
```

This retrieves the actual video data from the specified stream using the `get_media` operation.

## 8. Create an S3 Key for the Video File

```
s3_key =
f'videos/client1/{datetime.datetime.now().strftime("%Y-%m-%d_%H-%M-%S")}.mp4'
```

This line constructs the file name (key) under which the video data will be stored in Amazon S3. The key will look something like `videos/client1/2024-09-19_14-35-50.mp4`, indicating that the video belongs to client1 and was created at the given timestamp.

## 9. Store the Video Data in Amazon S3

```
s3_client.put_object(
    Bucket='client1-video-storage',
    Key=s3_key,
    Body=video_data['Payload'])
```

This line uploads the video data retrieved from Kinesis to an Amazon S3 bucket for persistent storage.

- **Function**: The `put_object` method of the S3 client uploads an object (in this case, the video data) to a specified S3 bucket.

## 10. Preprocessing

Preprocessing is essential for preparing raw video frames or images before sending them to Amazon Rekognition. The goal is to improve the quality and structure of the input so that the object detection model works optimally. Once the Lambda function is triggered, it downloads the video from S3, extracts frames from

the video, resizes them, compresses them, and reduces noise. These steps ensure the video frames are optimized before being sent to Rekognition.

1. **Frame Extraction from Video**
   - Why: Rekognition works on images or video snippets, but it cannot process continuous video streams directly. You need to extract individual frames or key segments from video files before sending them for analysis.
   - Use OpenCV in the Lambda functions to extract frames from the video file.
   - Extract frames at regular intervals (one frame per second or at key moments) to avoid sending too many redundant frames.
   - Adjust the interval based on the video content. If the video is fast-paced (with moving objects), you might want to extract more frames; for static environments, fewer frames might be sufficient.
2. **Video Compression and Resizing**
   - Why: Large video files can increase the cost and latency of sending data to Rekognition. Compressing and resizing frames can reduce the payload size and make processing faster.
   - Use OpenCV to resize video frames to a lower resolution (720p or even lower if high resolution isn't needed).
   - Compress the frames to a more efficient format (JPEG) with a balance between file size and quality.
3. **Noise Reduction and Frame Cleaning**
   - Why: If the video is captured in environments with poor lighting or noise (shadows, glare, etc.), applying noise reduction techniques can improve the quality of object detection.
   - Use Gaussian blur to reduce noise in frames.
   - Enhance contrast or brightness if needed to ensure objects are clearly visible.

## 11. Video Analysis

Amazon Rekognition is used to perform object detection on stored video files. This service can identify objects such as people, vehicles, and other items in the video, along with the timestamps of when they appear. Note: we are currently working on integrating client-specific custom object-detection using Amazon Sagemaker Object Detection using the client's stored video feed history. At this stage, we are using Ground Truth, an Amazon managed data labeling service where we use the Ground Truth web-based interface to draw bounding boxes and assign labels. Once sufficient labels have been assigned, Ground Truth uses automated labeling predictions.

**Why Amazon Rekognition?**

- **Pre-trained Models**: Rekognition has pre-trained models for detecting common objects, reducing the need for custom model development.
- **Scalability**: Rekognition can analyze video data at scale without the need to manage infrastructure.
- **Integration**: Easily integrated with S3 and Lambda for automated workflows.

Every time a new video is stored in S3, a Lambda function will trigger to start preprocessing and object detection using Rekognition.

1. `rekognition = boto3.client('rekognition')` creates a Rekognition client to interact with the Amazon Rekognition service.
2. `def analyze_video(video_s3_key)` function is designed to perform video analysis. It accepts a single argument, `video_s3_key`, which represents the key (or file name) of the video in the S3 bucket that needs to be analyzed.

3. This method initiates label detection on the video stored in the specified S3 bucket. Label detection refers to the identification of objects, scenes, and activities within the video. It returns a JobId, which uniquely identifies the video analysis job. This JobId can be used later to retrieve the results of the analysis or to monitor the job's status.

```
response = rekognition.start_label_detection(
     Video={'S3Object': {'Bucket': 'client1-video-storage', 'Name':
     video_s3_key}},
     NotificationChannel={
             'RoleArn': 'arn:aws:iam::account-id:role/RekognitionS3Role',
             'SNSTopicArn':
     'arn:aws:sns:region:account-id:rekognition-topic'})
return response['JobId']
```

- **Parameters**:
  - ○ `Video={'S3Object': {'Bucket': 'client1-video-storage', 'Name': video_s3_key}}` specifies the location of the video in the S3 bucket.
  - ○ `NotificationChannel={...}`: Rekognition can send notifications about the status of the label detection job (like when it's complete or if it fails).
    - ■ `RoleArn`: This is the IAM Role that grants Rekognition permission to access the S3 bucket and read the video file.
    - ■ `SNSTopicArn`: The Amazon SNS Topic ARN that Rekognition uses to send notifications regarding the progress of the video analysis. SNS (Simple Notification Service) is a messaging service that enables you to receive alerts or messages related to the video analysis process.

## 11. Post-Processing

Post-processing takes the output from Amazon Rekognition and refines it for better downstream use. This can include filtering out irrelevant detections, adding custom logic, and managing results in a meaningful way.

1. **Confidence Score Filtering**
   - Why: Amazon Rekognition returns multiple detected objects along with their confidence scores. To ensure the results are meaningful, you should filter out detections below a certain confidence threshold.
   - Set a confidence score threshold (e.g., 80%). Only process objects detected with confidence levels above this threshold.
2. **Result Aggregation and Summarization**
   - Why: If multiple frames are analyzed, you might get redundant or overlapping results. Aggregating the results can reduce clutter and make it easier for clients to review the data.
   - Group similar detections across consecutive frames. If the same object appears in multiple frames, aggregate the results to avoid duplication.
   - Create summaries like "Detected 3 people over the last 10 minutes" instead of listing each detection separately.

## 12. Storing Detected Metadata

Once Amazon Rekognition analyzes a video and detects objects, the system stores the analysis results along with additional metadata such as video location, camera ID, detected objects, and timestamps in Amazon DynamoDB. This ensures that the data is readily available for fast querying via the chatbot or other interfaces. To optimize querying by location and camera ID, we also create a Global Secondary Index (GSI).

**Steps for Storing Analysis Results:**

### 1. DynamoDB Resource Creation
- **`dynamodb = boto3.resource('dynamodb')`**: This creates a DynamoDB resource object that allows interaction with DynamoDB in a high-level, object-oriented manner compared to using a low-level client.
- **`table = dynamodb.Table('VideoAnalysis')`**: This creates a reference to the DynamoDB table named `VideoAnalysis`, where the metadata and analysis results will be stored.

### 2. Global Secondary Index (GSI) Creation for Efficient Querying
To efficiently query data by location and camera ID, we create a GSI on the `location` and `camera_id` attributes. This allows the system to efficiently retrieve data based on these attributes without scanning the entire table.

- **Partition Key**: `location`
- **Sort Key**: `camera_id`
- **Projection Type**: Set to ALL so that all attributes in the table can be retrieved through the GSI.

3. **Additional Attributes**
   The system captures and stores additional metadata in DynamoDB such as the timestamp and the list of objects detected by Rekognition in the video.

**Example `store_analysis_results` Function:**

```
import boto3
from datetime import datetime

# Initialize DynamoDB resource and table reference
dynamodb = boto3.resource('dynamodb')
table = dynamodb.Table('VideoAnalysis')

def store_analysis_results(job_id, video_key, analysis_results, location, camera_id):
    """
    Stores the analysis results and additional metadata in DynamoDB.

    Parameters:
    - job_id: The job ID associated with the video analysis task.
    - video_key: The key identifying the video file in S3.
    - analysis_results: The results from Rekognition (detected objects, labels, etc.).
    - location: The physical or logical location where the video was recorded.
    - camera_id: The identifier of the camera that recorded the video.
```

```
"""
# Generate the current timestamp
timestamp = datetime.now().strftime('%Y-%m-%d %H:%M:%S')

# Insert the analysis results and metadata into DynamoDB
table.put_item(
    Item={
        'job_id': job_id,
        'video_key': video_key,
        'detected_objects': analysis_results,  # Detected objects from Rekognition
        'location': location,  # Video recording location (GSI partition key)
        'camera_id': camera_id,  # Camera used (GSI sort key)
        'timestamp': timestamp  # Time of analysis
    }
)
```

## 13. MLOps

When using Amazon Rekognition for video analysis, the key challenge is ensuring that the system is automated, scalable, and easy to maintain. Even though we are using a pre-built model like Rekognition (meaning model versioning is not necessary), it still benefits from a robust MLOps pipeline for managing infrastructure changes, tracking events, automating deployments, and sending client-specific alerts.

Even though Model Versioning is not necessary when using Amazon Rekognition, you still need to version other components like the infrastructure, code, and configuration settings. Version control and CI/CD are critical for these reasons:

- Tracking Configuration Changes: Any changes to the S3 event triggers, object detection thresholds, or Lambda function logic need to be version-controlled and tested.
- Ensuring Reproducibility: If you need to roll back to a previous version of the pipeline (due to a bug or error), having a versioned history of code and configurations is essential.

    **1. Set Up a Code Repository**:
    - Use AWS CodeCommit to store the Lambda function code, infrastructure configuration files, and any other relevant resources.

    **2. AWS CodePipeline for CI/CD**
    - To manage continuous integration (CI) and continuous deployment (CD) of changes to the pipeline, specifically for managing changes in:
        - AWS Lambda functions (updates to preprocessing logic, object filtering, or post-processing results).
        - Infrastructure configurations (updates to S3 event triggers, IAM roles, DynamoDB schema for storing results).
    - **Why CodePipeline**:
        - Automated Testing and Deployment: Anytime someone updates any part of the code CodePipeline automatically runs tests and validates changes using AWS CodeBuild then deploys the updated code to our production environment by pushing the updated code to AWS Lambda. Whenever changes are committed to the repository, the pipeline is triggered, and the updated Lambda function is deployed. This reduces manual

intervention, ensures updates are thoroughly tested, and minimizes the risk of introducing errors.
- ● Tracking Changes: CodePipeline integrates with CodeCommit version control so each code change is tracked, and we can easily roll back to previous versions if needed.

### 3. Amazon SNS for Client-Specific Custom Alerts

- ● SNS is implemented to send real-time notifications when something significant is detected by Rekognition (like a person or vehicle is detected in a restricted area), or when there's a system error.
  - ● Real-Time Alerts: SNS allows you to notify clients immediately when their video footage detects certain events via email or SMS.
    - ● Requires a new SNS topic for each client or event type to be created in the SNS Console and client subscription to the topics.

### 4. Amazon CloudWatch for Monitoring and Logging

- ● Used for monitoring the performance, health, and usage of the entire event-driven pipeline. CloudWatch provides insights into:
  - ○ **API performance** (latency, error rates for Rekognition API calls).
  - ○ **Lambda function execution** (errors, timeouts, and resource utilization).
  - ○ **Custom metrics** (tracking how often a specific object is detected).
- ● **Why CloudWatch**:
  - ● **Monitoring**: CloudWatch allows you to monitor the health of your system in real-time and provides dashboards for visualizing important metrics. You can create custom metrics to track specific events, such as the number of objects detected by Rekognition, the processing time per video, or the latency of API calls.
  - ● **Logging**: CloudWatch also stores logs from Lambda functions and other services, helping you troubleshoot issues (Lambda errors or performance bottlenecks). By default, Lambda logs are sent to CloudWatch Logs, where you can review execution details.
  - ● **Alarms**: CloudWatch Alarms can be set to trigger alerts (via SNS) when certain conditions are met (if Rekognition latency exceeds a threshold or if a Lambda function fails).

## 14. Security and Access Control

Security and access control are paramount to ensure that clients only access their own data and that the system complies with security best practices.

### 1. AWS Cognito for Client Authentication and Authorization

- ● AWS Cognito is used to securely manage user sign-up and sign-in, as well as to provide access control for resources. It handles user authentication, ensuring that only authorized users can access the system.
- ● Once authenticated, Cognito generates a JWT token, which represents the user's identity and contains user details (claims like sub). This token is passed to the backend via API Gateway to validate their identity by extracting the `user_id` from the token and verifying that the client is authorized to access only their own data.
- ● Cognito supports multi-factor authentication, password policies, and integration with social identity providers (Google, Facebook) or enterprise directories (SAML).
- ● After authentication, Cognito works with IAM policies to control which resources a client can access.

## 2. IAM Roles and IAM Policies

- IAM Roles allow services like Lambda, API Gateway, and Cognito to securely access AWS resources such as S3, Rekognition, DynamoDB, and more.
- Each service in our pipeline has an assigned IAM Role, and each role has a set of permissions that dictate what actions the service can perform on specific resources. For example, the Lambda function that processes video clips has a role that grants it read access to S3, query access to DynamoDB, and invoke Rekognition API calls.
- IAM Policies Fine-Grained Access gives each client restricted access based on their identity, enforced by resource-level permissions and condition keys. For example, a client can only access their own video data in the S3 bucket by including the `user_id` in the S3 object key and writing an S3 policy that allows access only to objects that match the user's ID. The client can only query records in the DynamoDB table where their `user_id` matches the `user_id` attribute in the table. This can be enforced by conditionally allowing access based on the Cognito identity.

## 3. Lambda Role Permissions

- The Lambda function that handles video preprocessing, querying, and URL generation will have an IAM role with policies that grant:
    - S3 access: For reading and writing video files.
    - DynamoDB access: For querying video analysis data.
    - Rekognition access: For calling the Rekognition API to analyze video frames.

Example IAM Role Permissions for Lambda (json):

```json
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "s3:GetObject",
        "dynamodb:Query",
        "rekognition:DetectLabels"
      ],
      "Resource": "*"
    }
  ]
}
```

## 4. Encryption at Rest (S3 SSE) and in Transit (HTTPS)

- Encryption at rest ensures that the video data stored in S3 is protected, even if unauthorized parties gain access to the underlying storage systems.
- S3 Server-Side Encryption (SSE) is used to automatically encrypt all video data stored in the S3 bucket using AWS-managed encryption keys. You can enable SSE-S3 by default for the S3 bucket where video files are stored.
- Encryption in transit protects data while it is being transferred between services, such as between Kinesis, Lambda, S3, and Rekognition.

- All data transfers between AWS services are encrypted using HTTPS (SSL/TLS), ensuring that video data is secure during transmission. When Lambda functions, API Gateway, or clients access S3 or Rekognition, the requests are made using HTTPS, preventing data from being intercepted.

## 5. Network Isolation and Security

- Placing sensitive resources (such as video streams, processing services, or inference endpoints) inside a Virtual Private Cloud (VPC) provides network-level isolation, ensuring that only authorized traffic can access these resources.
- Control access to resources through security groups, network access control lists (NACLs), and VPC endpoints.
- To allow Lambda functions within the VPC to securely access S3 and DynamoDB, you can configure VPC endpoints. This allows traffic between Lambda and other AWS services to stay within the AWS network, without traversing the public internet.

# 15. Querying via Chatbot

Amazon Lex is used to build a chatbot interface that clients can use to query their video footage. Lex supports natural language processing, allowing users to search for specific objects (like "Show me videos with vehicles from last Tuesday").

- **Natural Language Processing**: Lex provides built-in NLP capabilities, enabling clients to query footage using natural language rather than SQL or complicated commands.
- **Integration with Lambda**: Lex integrates with Lambda for custom query execution. The Lambda function that Lex calls queries the DynamoDB table for metadata matching the client's search query. It will then return secure URLs to access the relevant video clips stored in S3.
- **Scalable and Secure**: Lex scales automatically and integrates with AWS Cognito for secure user access.

**Authentication and Authorization of the User**
Before querying the video data and generating presigned URLs, the system must first authenticate the user and ensure that they are authorized to access only their own video feed.

- **User Authentication**: The client's identity is verified using Amazon Cognito. This ensures the user is who they claim to be.
- **User Authorization**: The system checks if the authenticated user has permission to access the specific video feed (belonging to their account).

**Process:**

**1. Extract User Identity**
- **user_id = event['requestContext']['authorizer']['claims']['sub']**:
  - The Lambda function retrieves the user's unique ID (sub from an Amazon Cognito JWT token) from the event object. This user_id is used to identify the authenticated user.
  - This ID is passed through API Gateway and a Lambda authorizer, ensuring that only authenticated users can trigger the Lambda function.

**2. Validate User Permissions**
- **check_user_permissions(user_id)**:
  - This step involves checking whether the authenticated user is authorized to access the specific video clips.
  - This function would query IAM policies to validate the user's permissions. It ensures that only the rightful owner of the video data can access it.
  - If the user is unauthorized, the function returns a failed response indicating that access is denied.

## 3. Extract the Lex Chatbot Slots

- **event['currentIntent']['slots']['Object'], event['currentIntent']['slots']['StartTime'], event['currentIntent']['slots']['EndTime']**:
  - These slots come from Amazon Lex, which captures user inputs (such as the object they are searching for and the time range for the video query). The values are extracted from the event object.
  - **object_type**: The type of object the user wants to search for in the video feed ("person", "vehicle").
  - **start_time and end_time**: The time range within which the user wants to find video clips.

## 4. Query DynamoDB for Relevant Video Clips

DynamoDB is a fully managed NoSQL database that stores key-value and document-based data. The primary structure is based on tables, which consist of items (similar to rows in SQL), each with a set of attributes (similar to columns).

Key points about how DynamoDB works:

1. Partition Key and Sort Key: Every table has a primary key, which can either be a partition key (single attribute) or a composite key (a combination of a partition key and sort key). These keys are used to uniquely identify items.
2. No Joins or Foreign Keys: DynamoDB is optimized for fast lookups based on keys. However, it does not support joins or foreign key constraints like relational databases. Therefore, it's designed to denormalize data (store related data together), which can reduce the need for complex data merging.
3. Indexes: You can create secondary indexes (Global Secondary Indexes or Local Secondary Indexes) to query the data using attributes other than the primary key. This gives more flexibility when querying, but it still doesn't provide native join functionality.
4. Horizontal Scalability: DynamoDB is designed for horizontal scalability, meaning it can scale out easily to handle massive amounts of traffic.

In DynamoDB, querying for video clips requires filtering based on the object type, the time range, and the user ID. This is done using the KeyConditionExpression. The KeyConditionExpression filters results based on one or more primary keys (either the partition key or sort key) or indexes that have been defined on the DynamoDB table.
  - **Key('detected_objects').eq(object_type)**: Filters video clips based on the detected object type.
  - **Key('timestamp').between(start_time, end_time)**: Filters clips based on the time range specified by the user.
  - **Key('user_id').eq(user_id)**: Ensures that the query only returns video clips that belong to the authenticated user. This protects users from accessing video feeds they do not own.

For filtering by location or camera ID, we use a Global Secondary Index (GSI), which allows you to efficiently query based on non-primary key attributes (like `location` and `camera_id`), using them as indexed keys.

## 5. Generate Presigned URLs for S3 Video Access

- **Presigned URLs**:
    - For each video clip returned by the DynamoDB query, a presigned URL is generated using the S3 object key (video_key), allowing the user to temporarily access the video file without requiring direct AWS credentials.
    - The presigned URL will expire after a specified time, ensuring time-limited access to the video files.
- **s3.generate_presigned_url()**:
    - **Params**:
        - 'Bucket': 'client1-video-storage': The S3 bucket storing the video.
        - 'Key': video_key: The key (path) to the specific video file in the bucket.
    - **ExpiresIn=3600**: The URL expires after 1 hour, after which it will no longer work.

## 6. Return URLs via Chatbot Response

- The return function returns the list of video URLs back to the user through the Lex chatbot interface, confirming that the request has been fulfilled.
- The URLs are returned as plain text in the message, allowing the user to click and view the video clips in their browser or download them.

**Example of code:**

```
import boto3
from boto3.dynamodb.conditions import Key

# Initialize resources
dynamodb = boto3.resource('dynamodb')
table = dynamodb.Table('VideoAnalysis')

s3 = boto3.client('s3')

def lambda_handler(event, context):
    # Step 1: Extract User Identity
    user_id = event['requestContext']['authorizer']['claims']['sub']  # Cognito user ID (or use another IDP claim)

    # Step 2: Validate User Permissions
    # Assume a function check_user_permissions exists to verify the user's access rights.
    if not check_user_permissions(user_id):
        return {
            "dialogAction": {
                "type": "Close",
                "fulfillmentState": "Failed",
                "message": {
                    "contentType": "PlainText",
                    "content": "You are not authorized to access these video clips."
                }
            }
        }


    # Step 3: Extract the Lex chatbot slots
    object_type = event['currentIntent']['slots']['Object']
```

```
start_time = event['currentIntent']['slots']['StartTime']
end_time = event['currentIntent']['slots']['EndTime']

# Step 4: Query DynamoDB for relevant video clips within user's scope
response = table.query(
    KeyConditionExpression=Key('detected_objects').eq(object_type) &
                Key('timestamp').between(start_time, end_time) &
                Key('user_id').eq(user_id)  # Ensure that user can access only their own data
)

# Step 5: Generate Presigned URLs for video access from S3
video_urls = []
for item in response['Items']:
    video_key = item['video_key']
    presigned_url = s3.generate_presigned_url('get_object',
        Params={'Bucket': 'client1-video-storage', 'Key': video_key},
        ExpiresIn=3600)  # URL expires in 1 hour
    video_urls.append(presigned_url)

# Step 6: Return URLs via chatbot response
return {
    "dialogAction": {
        "type": "Close",
        "fulfillmentState": "Fulfilled",
        "message": {
            "contentType": "PlainText",
            "content": f"Here are your video clips: {', '.join(video_urls)}"
        }
    }
}
```

## 16. Monitoring and Logging

**Monitoring and logging** are essential aspects of managing the health, performance, and stability of our system, especially when dealing with serverless architectures and machine learning workloads.

- **Amazon CloudWatch Logs** is a service that collects and monitors log data from AWS services. It allows you to capture logs, filter log events, and search for errors or performance issues in your serverless infrastructure.
- Logs help capture errors and exceptions that occur during the execution of Lambda functions. This is essential for troubleshooting and diagnosing issues in real time. Logs can also track performance metrics such as execution time, memory usage, and API call latency, which can be used to identify bottlenecks or optimize the system. We use JSON format structured logging for easier parsing, filtering, and log analysis in CloudWatch Log Insights.
- **CloudWatch Alarms** allow you to track critical metrics and set alarms that notify you when a specific threshold is breached (like high error rates, long inference times, or throughput issues). When an alarm is triggered, it will automatically notify our team, take action (such as scaling up resources), or trigger other workflows. These are set up programmatically using boto3 SDK.
- **Key Metrics to Monitor:** Lambda function errors, model inference latency, Kinesis Stream throughput, and CPU/memory utilization.