

An Adaptive File Connector for Spark to Hunt for Cyber Attacks



Wojciech Indyk
Big Data Engineer, Hunters



Ada Sharoni
Senior Software Engineer
and Team Leader, Hunters

Ada Sharoni

Software Engineer Team Lead @ Hunters.ai

1. ML & Big Data
2. 7 Years, different security solutions:
 - a. Enterprise Network Security
 - b. WAF (Web Application Firewall)
 - c. Fraud Detection
3. Fun Fact: I started out as a Hardware Engineer

<https://twitter.com/AdaSharoni>

<https://www.linkedin.com/in/ada-sharoni-47ba26b8/>



Wojciech Indyk

Big Data Engineer @ Hunters.ai

- Big Data and Data Science for over 10 years
- 11 scientific papers regarding distributed machine learning systems
- I love open-source! ❤️

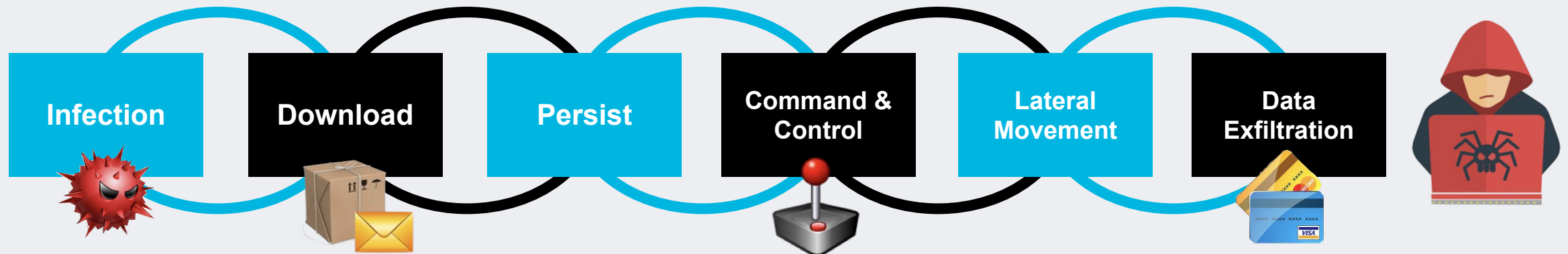
<https://pl.linkedin.com/in/wojciechindyk>



Hunters

Security Operations Platform

- Help security teams understand the full attack story
- Correlate existing telemetry and sources across surfaces
 - Network
 - Cloud
 - Email
 - Endpoint
 - SasS
 - etc



Variety of Data Sources

The screenshot shows the HuntersAI interface with a 'Stories' view. The table displays a sequence of events with associated data points:

Event	WHO	WHAT
+18 DNS Query to low reputation domain #crowdstrike-raw-events 1 lead		Requested Domain: hunters-ss0.xyz CS Agent ID: 2a92...55e7
+56 Azure sign-in marked as risky by Microsoft #azure-signin 1 lead	Properties User Principal Name: Larrymiller@hunterlab.onmicrc	Properties Risk Detail: userPassedMFADrivenByRiskBa Properties App Display Name: Azure Portal
+50 Azure disk snapshot download link generated #azure-activity 1 lead	Person Name: Larry Miller	Azure Appid: c44b4083-3bb0-49c1-b47d-974
+25 Azure VM extension execution #azure-activity 1 lead	Person Name: Larry Miller	Azure Appid: c44b4083-3bb0-49c1-b47d-974

The 'Add Data Flows' dialog box is shown, listing various products for selection:

- 1 PRODUCT
- 2
- 3

The list of products includes:

- ADP
- Cisco Umbrella
- Cloudflare
- CrowdStrike
- Elastic Beats
- Fortinet
- GitHub

Buttons for 'TEST CONNECTION' and 'SUBMIT' are visible at the bottom of the dialog.

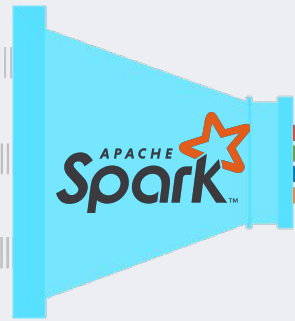
Streaming Security Data in Real-Time

Data Sources



Flexible Ingestion

- multiple formats
- multiple sources
- Streaming in real-time



ingestion

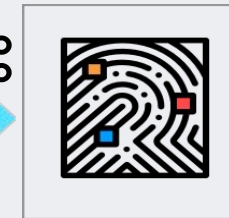
Data Lake



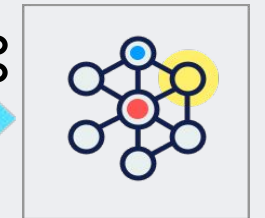
Detection Layer



Auto Investigation

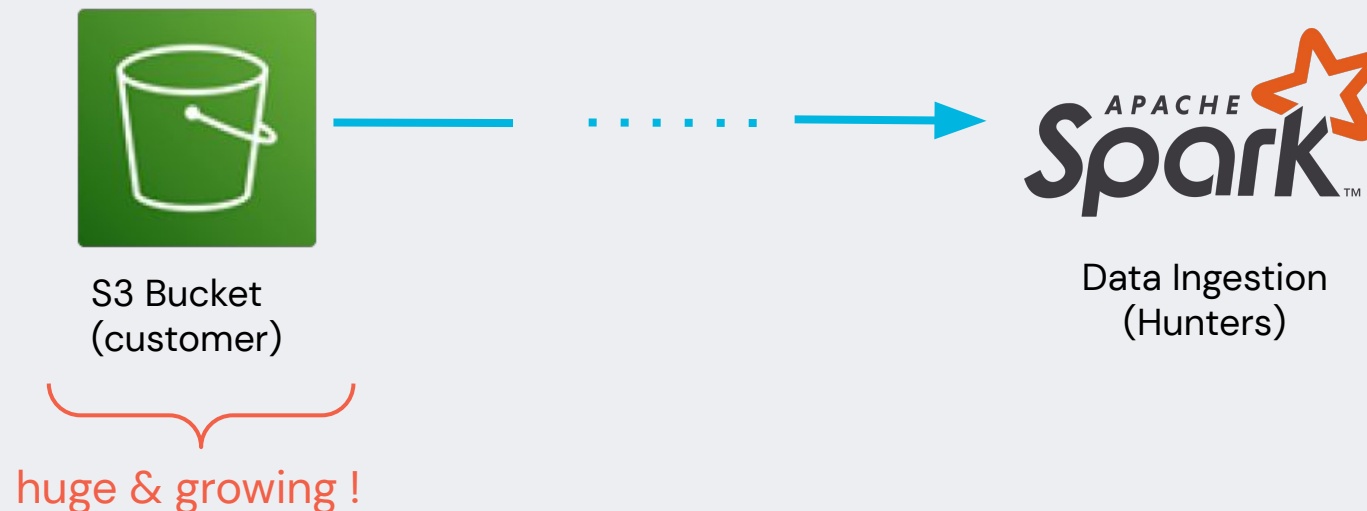


Knowledge Graph



Problem definition

1. 50% Of security data logs get shipped to object store directly (e.g. S3)
2. S3 Bucket belong to customers → no retention, huge bucket size
3. We need to detect live cyber attacks in real-time → streaming



Unfortunately

...

S3 is not HDFS !

1. S3 is an object store
2. Listing huge buckets is problematic:
 - a. Takes ~24h
 - b. Operation costs money

There Are Existing Solutions

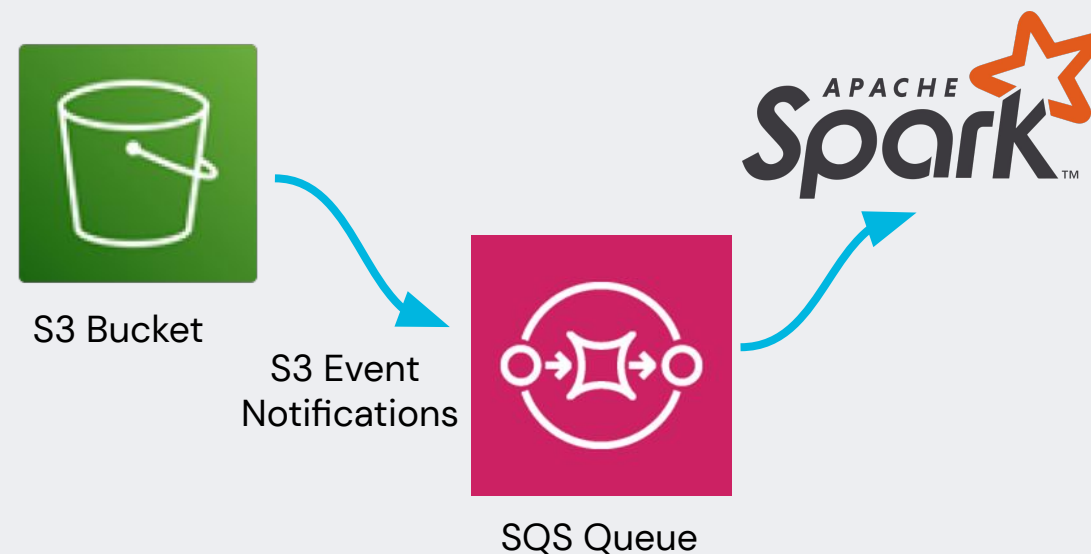
S3 Streaming via SQS

Open source: [qubole/s3-sqs-connector](https://github.com/qubole/s3-sqs-connector)

Managed: [Databricks's auto-loader](#)

Pain Points:

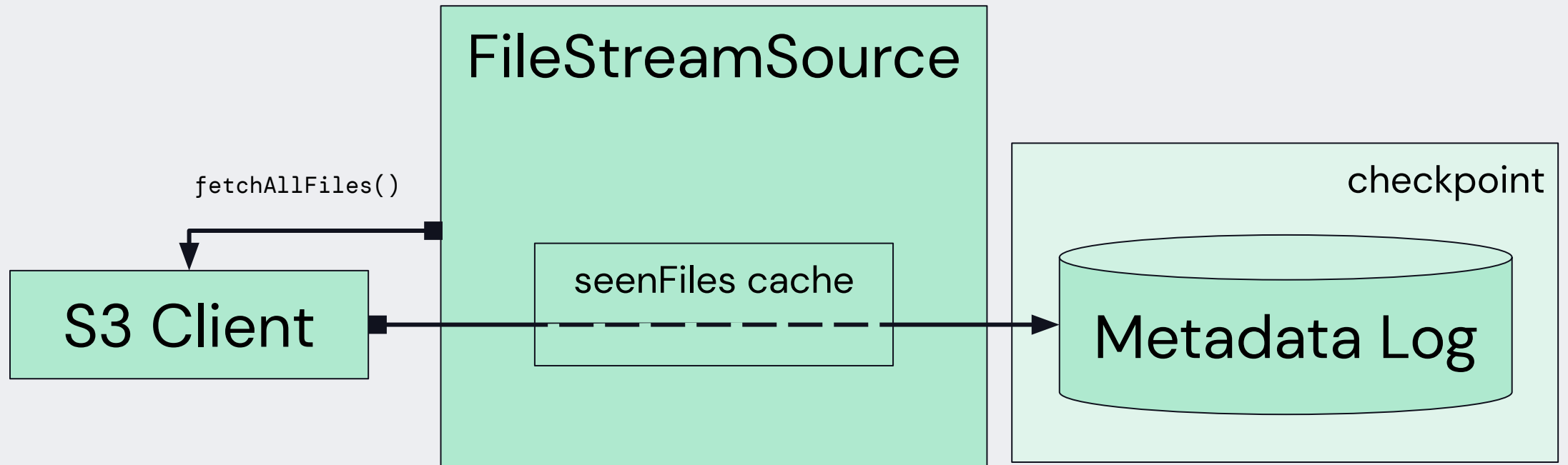
- Increases complexity of the system
- Customers don't want to grant sensitive permissions



So how can I
stream data from
s3 buckets ?

Spark-core current state

FileStreamSource - how does it work?

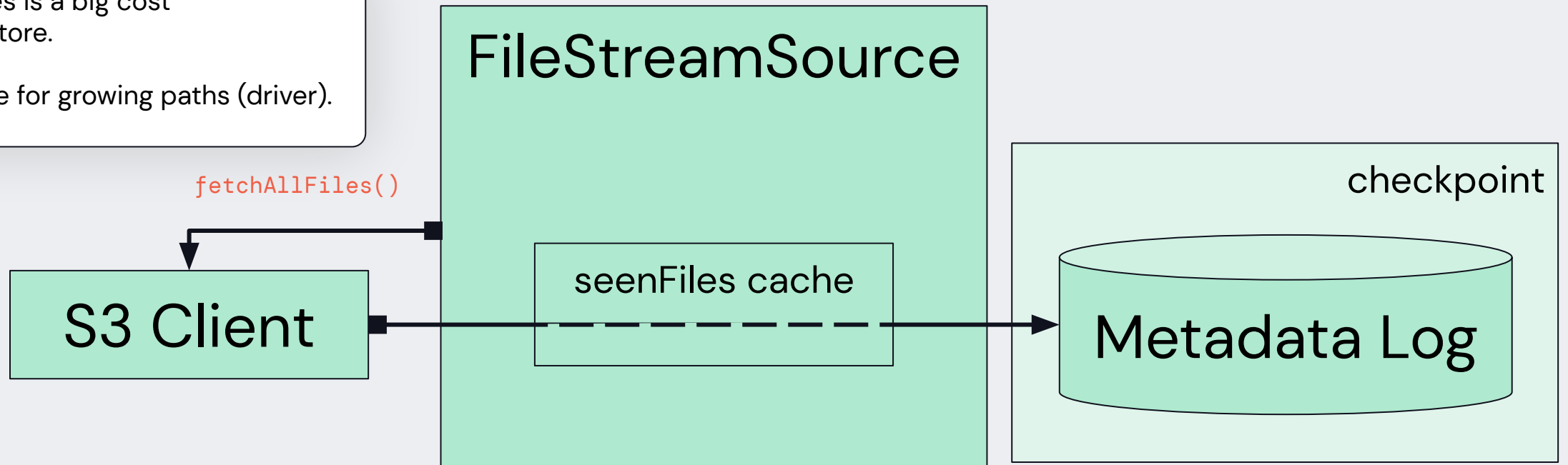


Spark-core current state

FileStreamSource - how does it work?

Fetch all files is a big cost
for object store.

Not scalable for growing paths (driver).



Solution design

Path structure

Directory structure for time-series data usually looks like:

s3:// <bucket_name> / <account>/ <service_name>/ <YYYY> / <MM> / <DD> /

Solution design

Time-series characteristic

Directory structure for time-series data usually looks like:

```
s3:// <bucket_name> / <account> / <service_name> / <YYYY> / <MM> / <DD> /
```



**this part is
dynamic in time !**

Solution design

Dynamic Path Generator

s3:// <bucket_name> / <account>/ <service_name>/ <YYYY> / <MM> / <DD> /



Static Part

(known upfront)



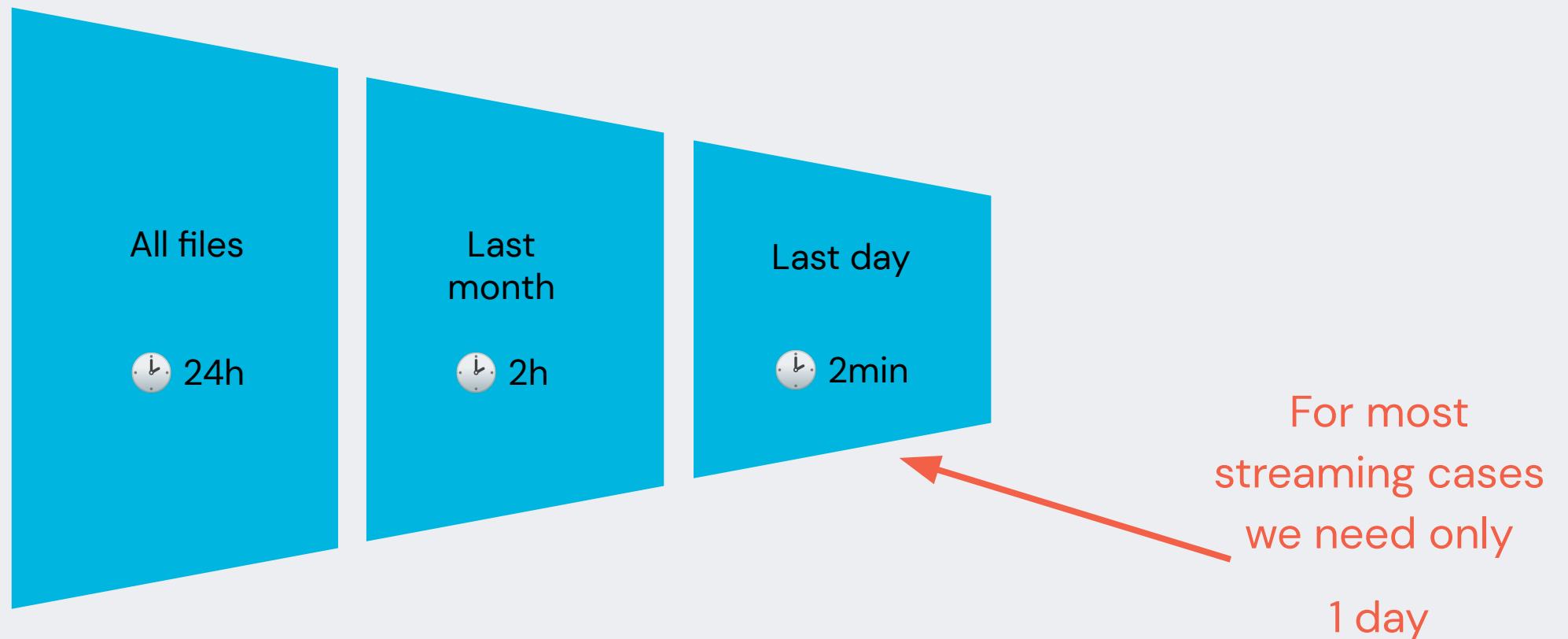
Dynamic Part

(includes time dimension)

Solution design

In streaming we only need to fetch the latest data

/ <YYYY> / <MM> / <DD> /



Recap

Recap

Solution design

If you:

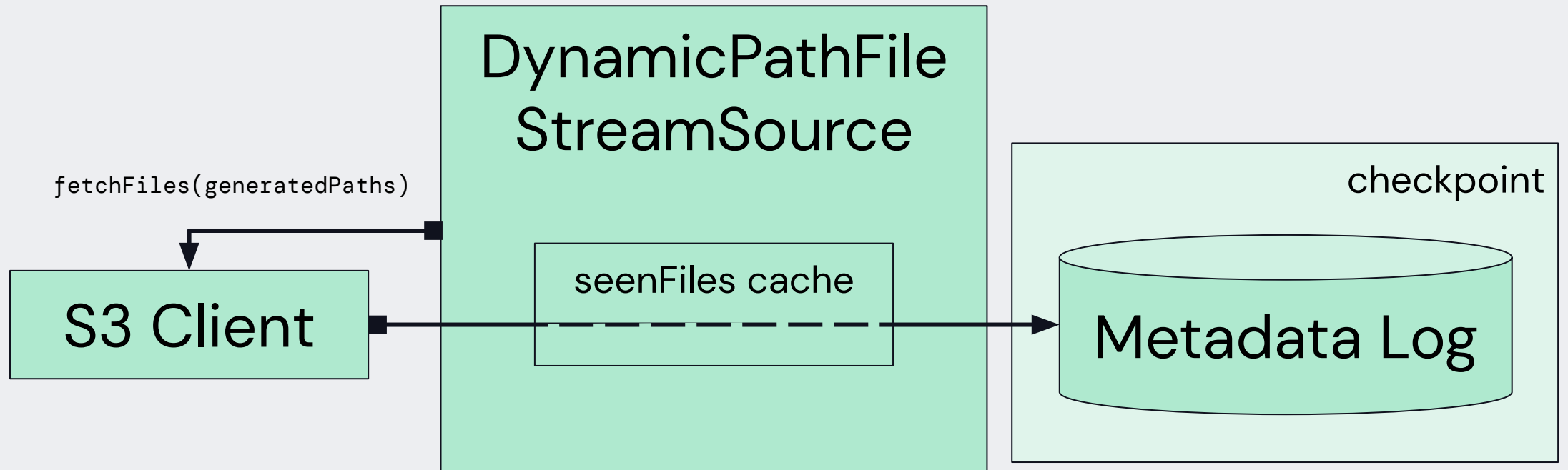
1. Are using Spark-Streaming
2. Know the structure of your S3 paths
3. Have dynamically changing part of your path
 - a. (And you know how it's changing, e.g. data partitioned by day)

Then you should pay attention to our adaptive file connector for spark !

Adaptive S3 connector for Spark

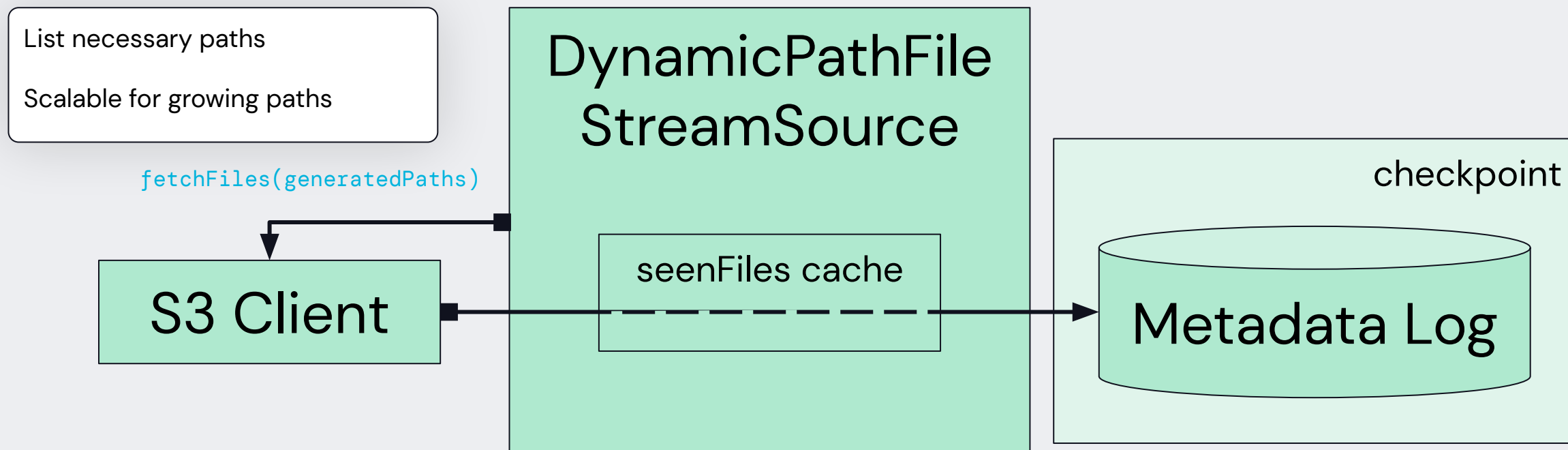
Adaptive S3 connector for Spark

Architecture



Adaptive S3 connector for Spark

Architecture




Adaptive S3 connector for Spark

Easy Usage

```
1  val readFileJsonStream = spark.readStream
2    .format("dynamic-paths-file")
3    .option("fileFormat", "json")
4    .schema(yourSchema)
5    .load("s3://my-bucket/prefix/{YYYY}/{MM}/{DD}")
```

Adaptive Source
Connector



```
7
8  val readFileJsonStream = spark.readStream
9    .format("json")
10   .schema(yourSchema)
    .load("s3://my-bucket/prefix/*/*/*")
```

Native Spark File
Source Connector



Adaptive S3 connector for Spark

Interface (Extending is also easy)

```
1  trait PathGenerator (  
2      def getPaths: Set[String]  
3  )  
4  
5  new DynamicPathGenerator extends PathGenerator (  
6      wildcardGlob: String,    // e.g. s3://my_bucket/data/{YYYY}/{MM}/{DD}/  
7      maxNumberOfDaysToRead: Int,  
8      timeZone: ZoneId        // timezone of data might be different from the server  
9  )  
10
```


Demo Time!

Demo - files to be processed

~/repos/hunters

12:50:39

```
aws --profile airflow-lab-to-hnt-test s3 --human-readable ls hunters-lab-audit/AWSLogs/591644025889/CloudTrail/us-west-1/2022/06/22/
```

Demo – Native Spark File Source

Spark Jobs (?)

User: root
Total Uptime: 7.4 h
Scheduling Mode: FIFO

▶ [Event Timeline](#)

Demo- Adaptive File Source

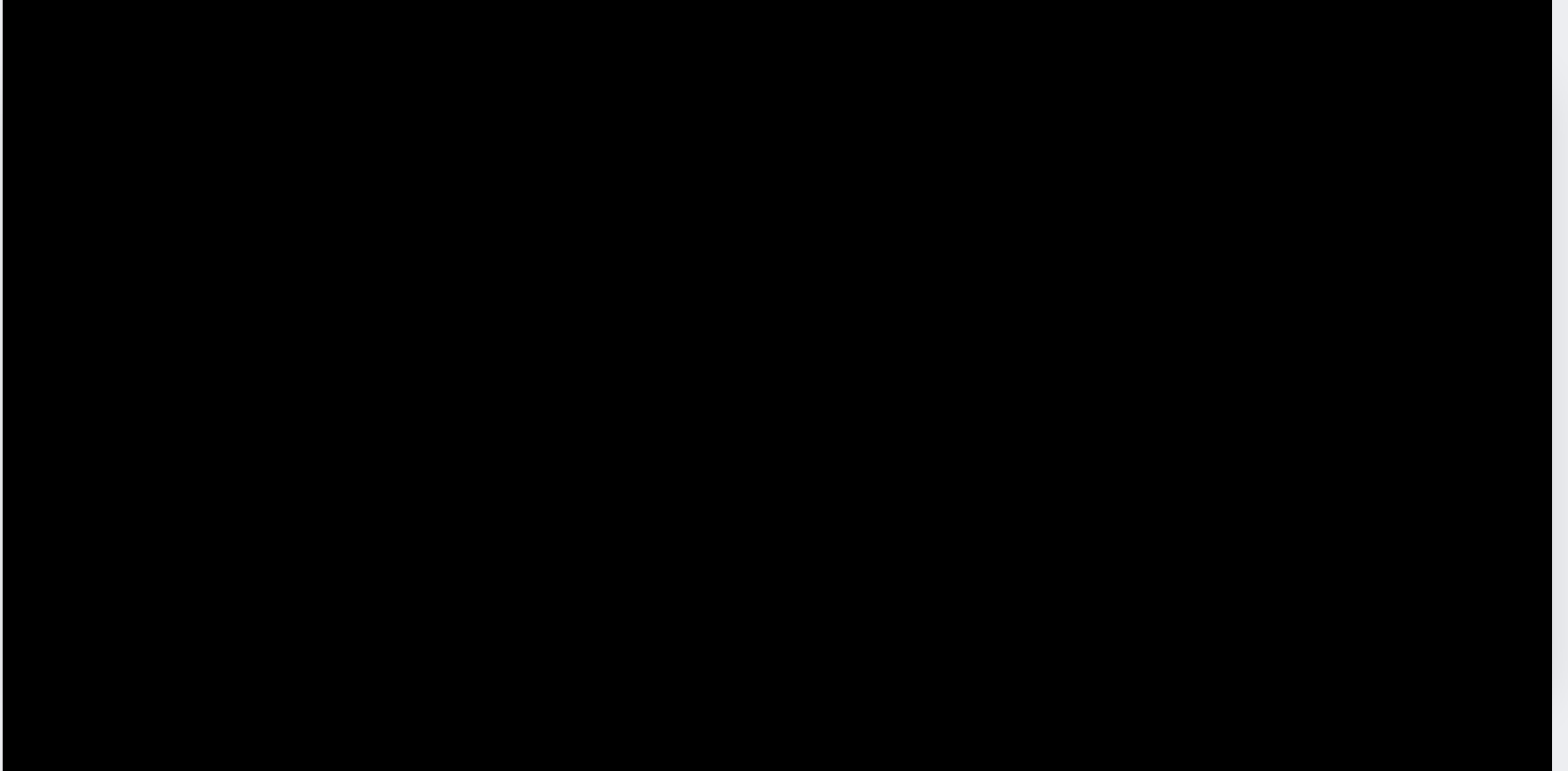
Streaming Query Statistics

Running batches for 25 minutes 52 seconds since 2022/06/22 10:26:04 (8 completed batches)

Name: spark
Id: f5b962e8-ef0b-46b7-9ba9-94297b821a42
Runid: 1251a54c-807b-4eeb-bc73-34dc2f465ba8



Simulated Real Time Attack



The solution is now Open-Sourced!

<https://github.com/hunters-ai/spark-adaptive-file-connector>

Sponsored by

 **HUNTERS**

Roadmap

Advanced S3 connector for Spark

1. Flexible specification of the time dimension- any spark-supported time formatting
2. Blacklist/whitelist pattern parameter
3. Extract more info from the bucket partitioning as extras columns - e.g. account, region, etc.
4. Try against Azure Blob Storage & Google Cloud Storage
5. Migrate Source from v1 to v2 for support of Continuous Processing mode

DATA+AI
SUMMIT 2022

Thank you



Wojciech Indyk

Big Data Engineer, Hunters

wojciech.indyk@hunters.ai



Ada Sharoni

Senior Software Engineer
and Team Leader, Hunters

ada@hunters.ai