## aws re: Invent

DECEMBER 2 - 6, 2024 | LAS VEGAS, NV

# Efficient incremental processing with Apache lceberg at Netflix

### Jun He

aws

(he/him) Staff Software Engineer Netflix

### Vaidy Krishnan

Strategic Account Manager AWS





# Big data

# Data accuracy

© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

# Data freshness

# **Cost efficiency**

# **Apache Iceberg**

# Netflix

© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

# **Netflix Maestro**

© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.



## About us



### Jun He

aws

(he/him) Staff Software Engineer Netflix



### Vaidy Krishnan

Strategic Account Manager AWS

## Outline

## 01 Introduction

- **02** Architectural design
- **03** Use cases and examples
- **04** Takeaways and future improvements



# Introduction



## Introduction

## Landscape of data insights at Netflix



## Data for business needs





## Streaming

aws

## Games





Ads





## **Common problems**

### Data accuracy



### Data freshness



## **Cost efficiency**



## Exabyte data warehouse

Business needs for new initiatives

More than \$150M per year



## **Big Data Analytics Platform**

## BDAP tech stack

![](_page_18_Picture_2.jpeg)

![](_page_18_Picture_3.jpeg)

Amazon EC2

![](_page_18_Picture_4.jpeg)

**Sec** snowflake

![](_page_18_Picture_5.jpeg)

![](_page_18_Picture_6.jpeg)

![](_page_18_Picture_7.jpeg)

![](_page_18_Picture_8.jpeg)

and other BDAP internal services

![](_page_18_Picture_10.jpeg)

Amazon S3

## **Existing solutions**

![](_page_19_Figure_1.jpeg)

## Incremental processing

Incremental processing is an approach to processing data in batches — but only new or changed data

- Capturing incremental data changes
- Tracking their states (i.e., whether a change is processed by a workflow or not)

## Incremental processing support (IPS)

## Key enablers

![](_page_21_Picture_2.jpeg)

![](_page_21_Picture_3.jpeg)

![](_page_21_Picture_4.jpeg)

aws

## Apache Iceberg

Iceberg is a high-performance format for huge analytic tables

- Open table format
- Open community standard
- Decentralized metadata
- Storage separation
- Schema and partition evolution
- Ease of data management

Rich information from the metadata layer

![](_page_22_Picture_10.jpeg)

## **Iceberg concepts**

- Catalog
- Table

- Snapshot
- Data file
- Partition

![](_page_23_Figure_6.jpeg)

## **Iceberg concepts**

- Catalog
- Table

- Snapshot
- Data file
- Partition

![](_page_24_Figure_6.jpeg)

## Iceberg concepts

- Catalog
- Table

- Snapshot
- Data file
- Partition

![](_page_25_Figure_6.jpeg)

aws

## What is Maestro

- Horizontally scalable workflow orchestrator
- Oversee the entire lifecycle of a workflow
- Serverless execution from user perspective
- Support both acyclic and cyclic workflows
- Multiple out-of-the-box reusable patterns
- Designed for expansion and integrations with others

![](_page_26_Picture_8.jpeg)

## Maestro

## Netflix data/ML workflow orchestrator

Alpha in September 2021
Beta in April 2022
Migrated all traffic in October 2023
GA in November 2022
GA in November 2022
Description
Open to the public in July 2024
Description
Desc

New features in progress

Started in March 2020

MAESTRO

### YAML definition

#### Description:

This is a sample description to [a Maestro workflow](https://github.com/Netflix/maestro)

#### cron: '@daily'

Workflow:

```
id: demo.pipeline
```

```
my_query: INSERT OVERWRITE TABLE
${TARGET TABLE} SELECT USER ID, DATE INT,
```

```
SUM(watch_time) FROM ${SOURCE_TABLE} WHERE
```

```
is_active=true AND DATE_INT > DATE_MINUS_7
```

group by USER\_ID, DATE\_INT;

```
jobs:
```

_	oh	
	00	

```
id: job.1
```

```
type: Spark
```

```
spark:
```

script: \${my\_query}

\$ scheduler push workflow.yaml
\$ scheduler run demo.pipeline my\_query="SELECT
1"

Workflow run has been successfully submitted. View on data.netflix: https://data.netflix.net/maestro/sandbox/demo.p

### YAML definition

Des	cription	<b>1:</b>

This is a sample description to [a Maestro

#### Trigger:

cron: '@daily'

```
id: demo.pipeline
my_query: INSERT OVERWRITE TABLE
${TARGET_TABLE} SELECT USER_ID, DATE_INT,
SUM(watch time) FROM ${SOURCE TABLE} WHERE
```

is\_active=true AND DATE\_INT > DATE\_MINUS\_7
group by USER ID, DATE INT;

jobs:

- job:
  - id: job.1
  - type: Spark
  - spark:
    - script: \${my\_query}

\$ scheduler push workflow.yaml
\$ scheduler run demo.pipeline my\_query="SELECT
1"

Workflow run has been successfully submitted. View on data.netflix: https://data.netflix.net/maestro/sandbox/demo.p ipeline/instances/1

### YAML definition

Description:		
This is a sample description to [a Maestro		
<pre>workflow](https://github.com/Netflix/maestro)</pre>		
Trigger:		
cron: '@daily'		
Workflow:		
id demo nineline		
my_query: INSERT OVERWRITE TABLE		
<pre>\${TARGET_TABLE} SELECT USER_ID, DATE_INT,</pre>		
SUM(watch_time) FROM \${SOURCE_TABLE} WHERE		
is_active=true AND DATE_INT > DATE_MINUS_7		
group by USER_ID, DATE_INT;		

- job:
  - id: job.1
  - type: Spark
  - spark:
    - script: \${my\_query}

\$ scheduler push workflow.yaml
\$ scheduler run demo.pipeline my\_query="SELECT
1"

Workflow run has been successfully submitted. View on data.netflix: https://data.netflix.net/maestro/sandbox/demo.p ipeline/instances/1

### YAML definition

#### Description:

This is a sample description to [a Maestro workflow](https://github.com/Netflix/maestro)

#### Trigger:

cron: '@daily'

#### Workflow:

```
id: demo.pipeline
my_query: INSERT OVERWRITE TABLE
${TARGET_TABLE} SELECT USER_ID, DATE_INT,
SUM(watch_time) FROM ${SOURCE_TABLE} WHERE
is_active=true AND DATE_INT > DATE_MINUS_7
group by USER_ID, DATE_INT;
jobs:
```

### - job:

```
id: job.1
```

```
type: Spark
```

```
spark:
```

```
script: ${my_query}
```

\$ scheduler push workflow.yaml
\$ scheduler run demo.pipeline my\_query="SELECT
1"

Workflow run has been successfully submitted. View on data.netflix: https://data.netflix.net/maestro/sandbox/demo.p

ipeline/instances/1

### YAML definition

#### Description:

This is a sample description to [a Maestro workflow](https://github.com/Netflix/maestro)

#### Trigger:

cron: '@daily'

#### Workflow:

```
id: demo.pipeline
my_query: INSERT OVERWRITE TABLE
${TARGET_TABLE} SELECT USER_ID, DATE_INT,
SUM(watch_time) FROM ${SOURCE_TABLE} WHERE
is_active=true AND DATE_INT > DATE_MINUS_7
group by USER_ID, DATE_INT;
jobs:
```

### - job:

id: job.1

```
type: Spark
```

spark:

script: \${my\_query}

\$ scheduler push workflow.yaml

\$ scheduler run demo.pipeline my\_query="SELECT
1"

Workflow run has been successfully submitted. View on data.netflix: https://data.netflix.net/maestro/sandbox/demo.p

ipeline/instances/1

### YAML definition

#### Description:

This is a sample description to [a Maestro workflow](https://github.com/Netflix/maestro)

#### Trigger:

cron: '@daily'

#### Workflow:

```
id: demo.pipeline
my_query: INSERT OVERWRITE TABLE
${TARGET_TABLE} SELECT USER_ID, DATE_INT,
SUM(watch_time) FROM ${SOURCE_TABLE} WHERE
is_active=true AND DATE_INT > DATE_MINUS_7
group by USER_ID, DATE_INT;
jobs:
```

```
- job:
```

```
id: job.1
```

```
type: Spark
```

```
spark:
```

```
script: ${my_query}
```

```
$ scheduler push workflow.yaml
$ scheduler run demo.pipeline my_query="SELECT
1"
```

Workflow run has been successfully submitted. View on data.netflix: https://data.netflix.net/maestro/sandbox/demo.p ipeline/instances/1

## Workflow orchestrator for everyone

![](_page_34_Figure_1.jpeg)

![](_page_34_Picture_2.jpeg)

### Interfaces

YAML Python JAVA REST GraphQL Maestro UI Integrations

Shell scripts Notebooks Java Docker SQL Dynamic

Parameterized workflows SEL Input/output Parameters

[]

### Extensible

Job templates Bring your own compute Async execution

![](_page_34_Picture_12.jpeg)

## **IPS Product briefs**

Provide a clean and easy-toadopt solution for the efficient incremental processing

- Data freshness
- Data accuracy
- Cost efficiency

![](_page_35_Picture_5.jpeg)
# Architectural design



# Design of incremental processing support (IPS)

Efficient incremental change capturing

- DO NOT read any real data from the table(s)
- Use Apache Iceberg metadata layer

Best user experiences

- Decoupling change capturing from business logics
- Clean interfaces over Maestro
- Language and engine agnostic





# Design of incremental change capturing

Apache Iceberg metadata provides all the needed information related to the new/updated data files

- Upper and lower bounds of a given column (partitioned or unpartitioned)
- Data file size, location

Build a zero-data-copy table including only change data

- No need for real data access
- No security concerns

# Incremental change capturing

Iceberg table snapshot Immutable data files Hidden partitions



# Incremental change capturing

New "append" snapshot Added immutable data files Mutated partitions



# Incremental change capturing

Empty table with the same schema

Add data file references to the change table

- Zero data copy
- No data access
- Same partition layout





### Alternatives

- Apache Iceberg + Flink
- Apache Iceberg + Spark Procedure (create\_changelog\_view)
- Apache Iceberg + Spark Structure Streaming

## Show me the code

Table source = catalog.loadTable(sourceTableIdentifier);

List<Snapshot> snapshots = StreamSupport.stream(

# Get a list of append snapshots

SnapshotUtil.ancestorsBetween(source, toSnapshotIdInclusive, fromSnapshotIdExclusive).spliterator(), false)

.filter(snapshot -> snapshot.operation().equals(DataOperations.APPEND))

.collect(Collectors.toList());

Table icdcTable = createChangeTable(catalog, source, icdcTableName));

AppendFiles af = icdcTable.newAppend();

for (Snapshot snapshot : snapshots) {

for (DataFile df : snapshot.addedDataFiles(source.io())) {
 af.appendFile(df);

# Add data files to the ICDC table

### af.commit();

## **Emerging patterns**

- Incrementally process the captured incremental change data and directly append it to the target table
- Use captured incremental change data as the row level filter list to remove unnecessary transformation

• Use the captured range parameters in the business logic



# Pattern 1: Incrementally append captured change data



# Pattern 1: Incrementally append captured change data



# Pattern 2: Use captured change data as row level filters



### Pattern 3: Use the captured range parameters



# **Onboarding concerns**

Compatibility with existing workflows

Extra costs

- Development cost
- Operational cost
- Maintenance cost

# Maestro incremental processing solution (IPS)

### Interfaces

- Table(s)
  - Table name(s) passed to user business logic over Maestro parameters
- IpCapture step type
  - Capture changes from last checkpoint
- IpCommit step type
  - Commit checkpoint for this run





### **Best user experiences**

Users love to have freedom to use their own approaches to solve their problems

- Fully compatible with existing workflows
- Don't ask to use specific languages
- Don't ask to choose a specific compute engine
- Don't ask to rewrite the business logic
- Don't ask to re-architecture the existing data pipelines
- Low code solution







### Workflow definition

#### Workflow:

id: auto.remediation.pipeline

#### iobs

- subworkflow:

id: normal\_workflow
workflow id: demo.pipeline

failure\_mode: IGNORE\_FAILURE

id: check\_status

type: NoOp

- subworkflow:
  - id: recovery

workflow\_id: example.recovery.wf

- subworkflow:
  - id: re\_run
  - workflow\_id: demo.pipeline
- job:
  - id: end

```
type: NoOp
```

```
dag:
```

end

- normal\_workflow -> check\_status
- check\_status:

```
IF params.getFromStep('normal_workflow',
'MAESTRO_STEP_STATUS') ==
'COMPLETED_WITH_ERROR': recovery -> re_run ->
```

#### OTHERWISE: end

### Workflow definition

### Workflow:

id: auto.remediation.pipeline

### jobs:

- subworkflow:
  - id: normal\_workflow
  - workflow\_id: demo.pipeline
  - failure\_mode: IGNORE\_FAILURE
- job:
  - id: check\_status
  - type: NoOp
- subworkflow:
  - id: recovery
  - workflow\_id: example.recovery.wf

### - subworkflow:

- id: re\_run
- workflow\_id: demo.pipeline
- job:
  - id: end

### dag:

- normal\_workflow -> check\_status
- check\_status:
- IF params.getFromStep('normal\_workflow',
  'MAESTRO\_STEP\_STATUS') ==
  'COMPLETED\_WITH\_ERROR': recovery -> re\_run ->
  end

#### OTHERWISE: end







© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

### Workflow definition

### Description:

This is a sample description to [a Maestro workflow] (https://github.com/Netflix/maestro)

### Trigger:

```
cron: '@daily'
```

### Workflow:

```
id: incremental.processing.pipeline
```

### jobs:

### - job:

```
id: capture_change
```

```
type: IpCapture
```

```
ip_capture:
```

```
source_tables:
```

- membership\_table

capture\_mode: ICDC

```
lata_operations:
```

### - subworkflow:

id: auto.remediation.pipeline
workflow\_id: auto.remediation.pipeline
SOURCE\_TABLE: icdc.membership\_table.01
my\_query: \$S3{my\_script.sql}

```
id: commit_checkpoint
    type: IpCommit
    ip_commit:
        capture_step: capture_change
dag: sequential
```



### Workflow definition

### Description:

This is a sample description to [a Maestro workflow] (https://github.com/Netflix/maestro)

### Trigger:

cron: '@daily'

### Workflow:

id: incremental.processing.pipeline

### jobs:

id: capture\_change

type: IpCapture

ip\_capture:

source\_tables:

- membership\_table

capture\_mode: ICDC

data\_operations:

- append

#### - SUDWORKILOW:

id: auto.remediation.pipeline
workflow\_id: auto.remediation.pipeline
SOURCE\_TABLE: icdc.membership\_table.01
my\_query: \$S3{my\_script.sql}

- job:

id: commit\_checkpoint
type: IpCommit
ip\_commit:
 capture\_step: capture\_change

dag: sequential

### Workflow definition

### Description:

This is a sample description to [a Maestro workflow] (https://github.com/Netflix/maestro)

### Trigger:

```
cron: '@daily'
```

### Workflow:

```
id: incremental.processing.pipeline
```

### jobs:

### - job:

```
id: capture_change
```

```
type: IpCapture
```

- ip\_capture:
  - source\_tables:

- membership\_table
- capture\_mode: ICDC
- data\_operations:
  - append
- subworkflow:

id: auto.remediation.pipeline
workflow\_id: auto.remediation.pipeline
SOURCE\_TABLE: icdc.membership\_table.01
my\_query: \$S3{my\_script.sql}

id: commit\_checkpoint

- type: IpCommit
- ip\_commit:
  - capture\_step: capture\_change



# **IPS highlights**

Efficient incremental change capturing

Handle late arriving data

Language and engine agnostic

Clean interface

aws

Compatible with existing workflows Low onboarding cost



# Use cases and examples



# Two-stage pipeline example



# Stage 1 workflow with IPS



Original SQL query in playback_daily_workflow:	New SQL query in ips_playback_daily_workflow:
Upost	Upgot
INSERT OVERWRITE	MERGE INTO
SELECT \${business_logic}	SELECT \${business_logic}
FROM playback_table	FROM playback_icdc_table WHEN MATCHED
	<pre>\${deduplication_logic}</pre>

Original SQL query in playback_daily_workflow:	New SQL query in ips_playback_daily_workflow:
Unset INSERT OVERWRITE	Unset MERGE INTO
TABLE playback_daily_table SELECT \${business_logic}	TABLE playback_daily_table SELECT \${business_logic}
playback_table	playback_icdc_table
	<pre>\${deduplication_logic}</pre>

Original SQL query in playback\_daily\_workflow:

Unset

INSERT OVERWRITE

TABLE playback\_daily\_table

SELECT

```
${business_logic}
```

FROM

aws

playback\_table

```
New SQL query in ips_playback_daily_workflow:

Unset

MERGE INTO

TABLE playback_daily_table

SELECT

${business_logic}

FROM
```

```
WHEN MATCHED ${deduplication_logic}
```

# Two-stage pipeline example



# Stage 2 workflow with IPS



Original SQL query:	New SQL query:
Lincot	Uncot
INSERT OVERWRITE	MERGE INTO
SELECT \${business_logic} FROM playback_daily_table	<pre>SELECT      \${business_logic} FROM      playback_daily_table JOIN      playback_daily_icdc_table ON      \${aggregation_group_by_keys} WHEN MATCHED      \${deduplicate_logic}</pre>
## Changes of write job on stage 2

## Original SQL query:

#### Unset

INSERT OVERWRITE TABLE playback\_daily\_agg\_table SELECT

```
${business_logic}
```

FROM

aws

```
playback_daily_table
```

```
New SQL query:
Unset
MERGE INTO
       TABLE playback_daily_agg_table
SELECT
       ${business_logic}
FROM
       playback daily table
JOIN
       playback_daily_icdc_table
ON
       ${aggregation_group_by_keys}
```

#### \${deduplicate\_logic}

## Changes of write job on stage 2

## Original SQL query:

#### Unset

INSERT OVERWRITE TABLE playback\_daily\_agg\_table SELECT

```
${business_logic}
```

FROM

```
playback_daily_table
```

```
New SQL query:
Unset
MERGE INTO
       TABLE playback_daily_agg_table
SELECT
       ${business_logic}
FROM
       playback_daily_table
JOIN
       playback_daily_icdc_table
ON
```

#### WHEN MATCHED

\${deduplicate\_logic}

## Multi-stage workflow with IPS



## Multi-table use case



## Multi-table use case



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

# Takeaways and future improvements



- Efficient capturing: Iceberg metadata enables incremental processing without accessing data
- Simplified adoption: Decoupling change capture reduces complexity
- Improved experience: Clean interfaces improve usability
- New patterns: IPS patterns apply to many scenarios



## **Future improvements**

Move from Table to View to reduce maintenance

Support other types of snapshots with the new version of Iceberg table specs

Work with Iceberg community to share the approach in Iceberg cookbooks

Multi-stage data backfill support



## Acknowledgment **Drew Goya** Andy Chu Charles Smith Eva Tse Jun He Kyoko Shimada Netix **Ely Spears Yingyi Zhang** ieam **Zhuoran Dong** Anjali Norwood Brittany Truong Abhinaya Shetty Prashanth Ramdas Bharath Mummadisetty

aws

## Thank you!

#### Jun He

aws

Staff Software Engineer Netflix in linkedin.com/in/jheua



#### Vaidy Krishnan

Strategic Account Manager AWS Please complete the session

survey in the mobile app

**ፚፚፚፚ** 

