



# Streaming analytics for time-series data

And hopefully we don't choke

John Mousa

Sr. Solutions Architect, Digital Native Businesses  
Amazon Web Services

# John Mousa



Engineer at heart, enabling  
customers in DACH and  
world wide to bridge the  
worlds of microservices and  
analytics



/in/johnmousa

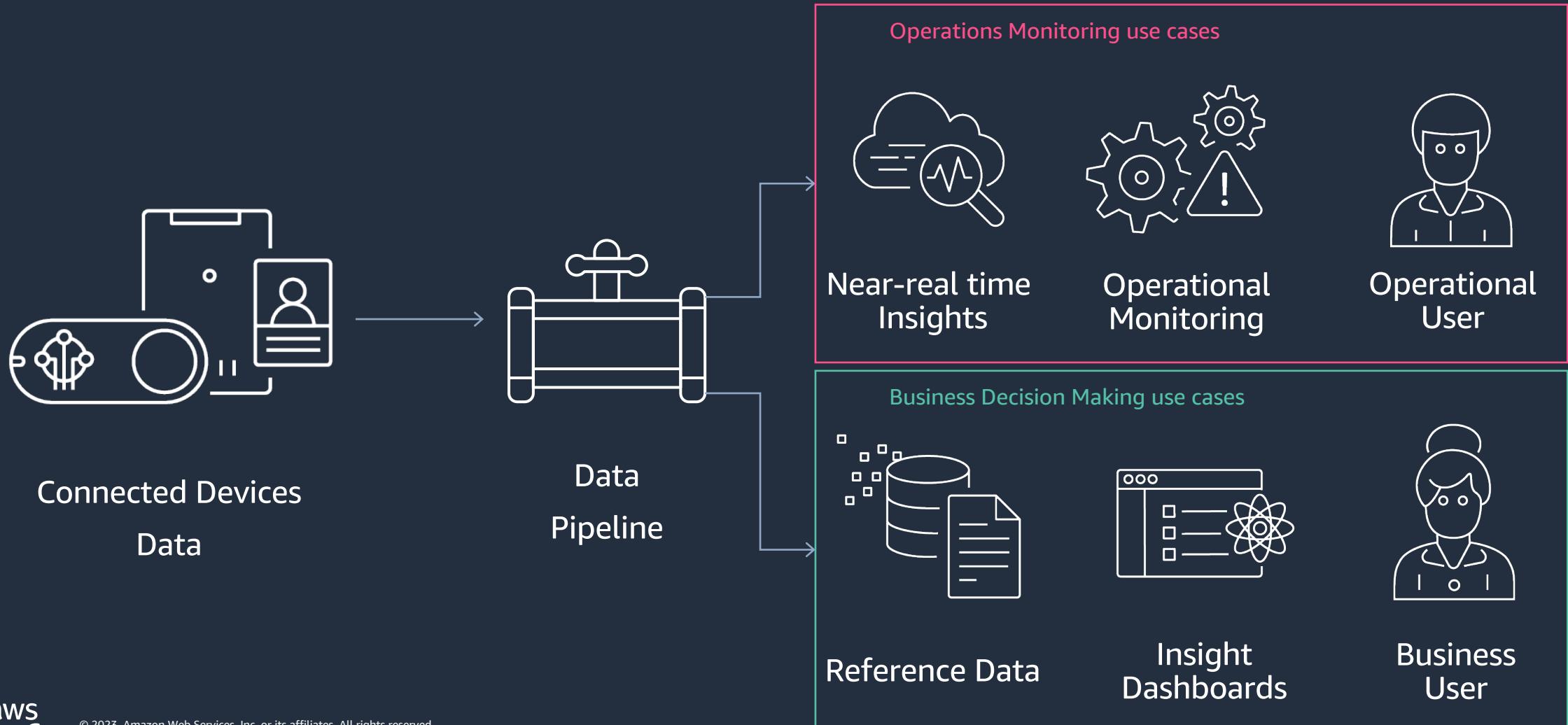


/johnmousa



© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

# Collecting telemetry data from smart meters



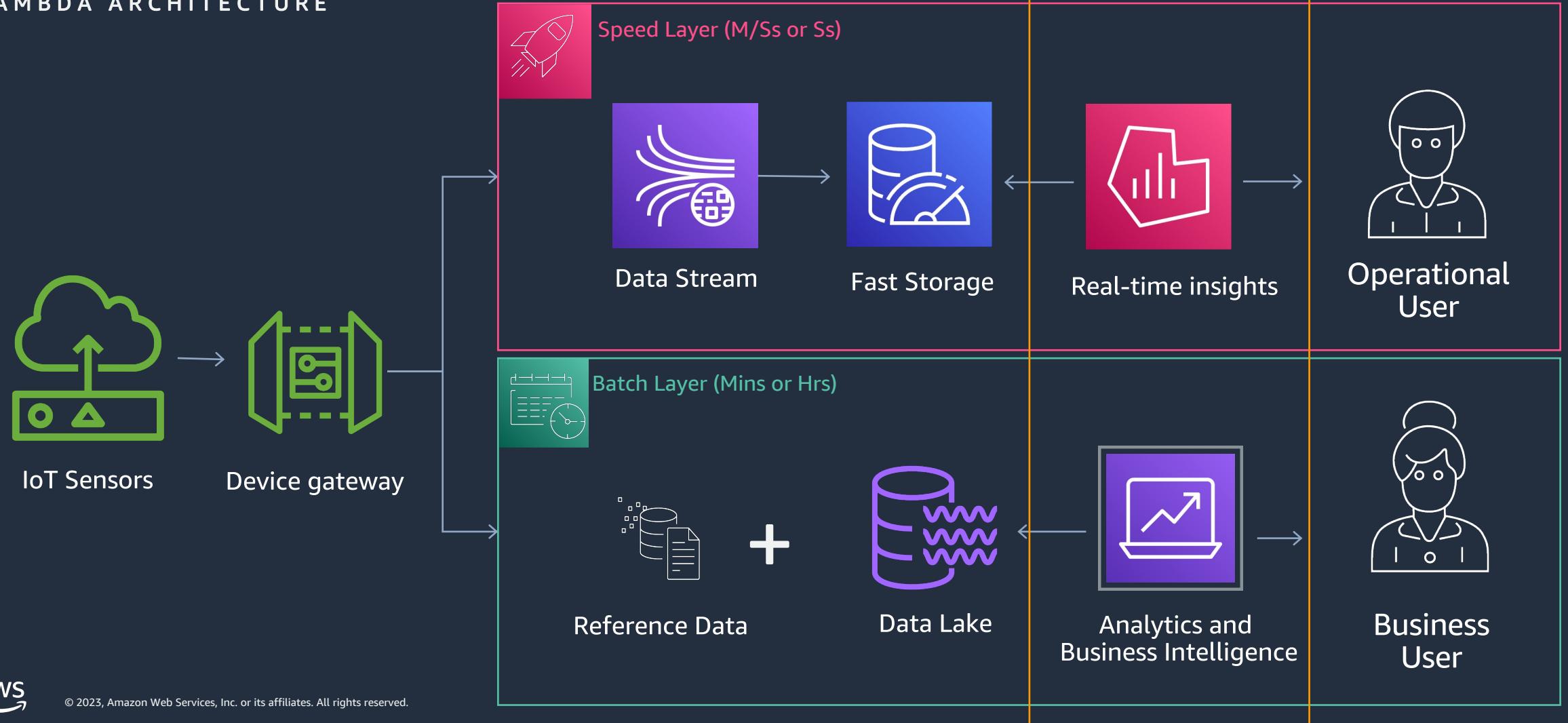
# Architecture



© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

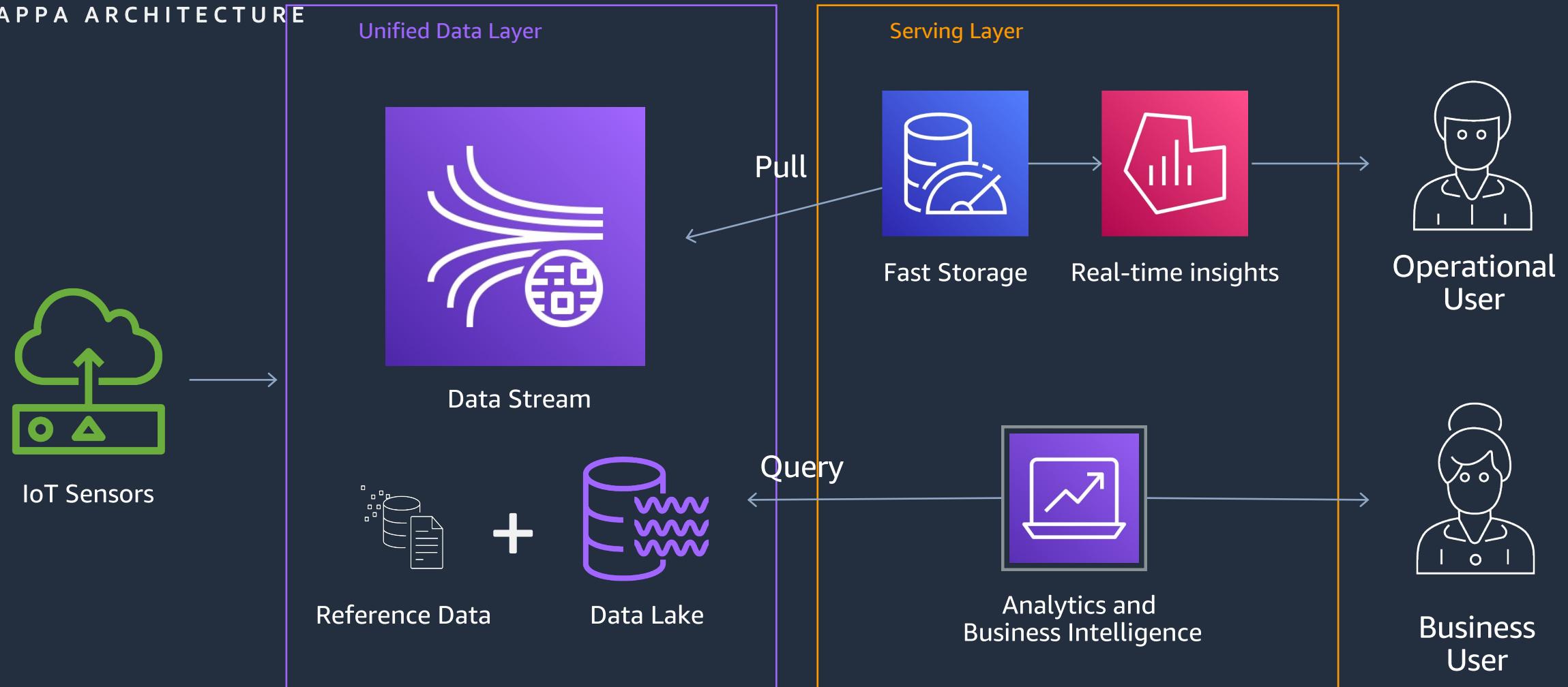
# Streaming data architectures

## LAMBDA ARCHITECTURE



# Streaming data architectures

## KAPPA ARCHITECTURE



# Streaming data architectures

$\lambda$

**Batch layer is built for big data**

**Result in overhead due to code  
and resource duplication**

$K$

**No redundancies**

**Batch workloads needs to be in a  
sink**

# Solutions



© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.



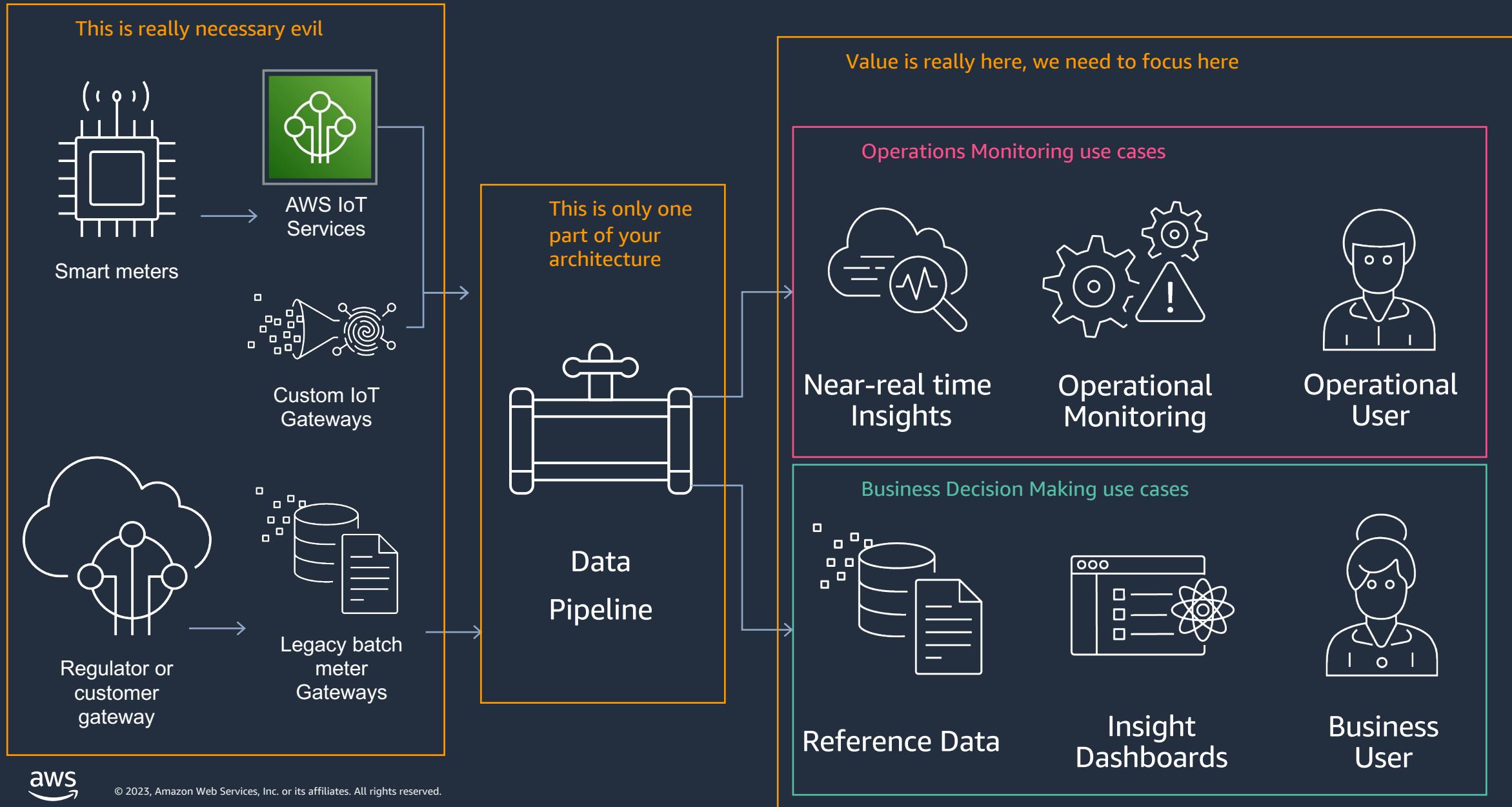
Streaming is Expensive

Streaming is Lossy

Hard to Operate,  
Scale and Maintain



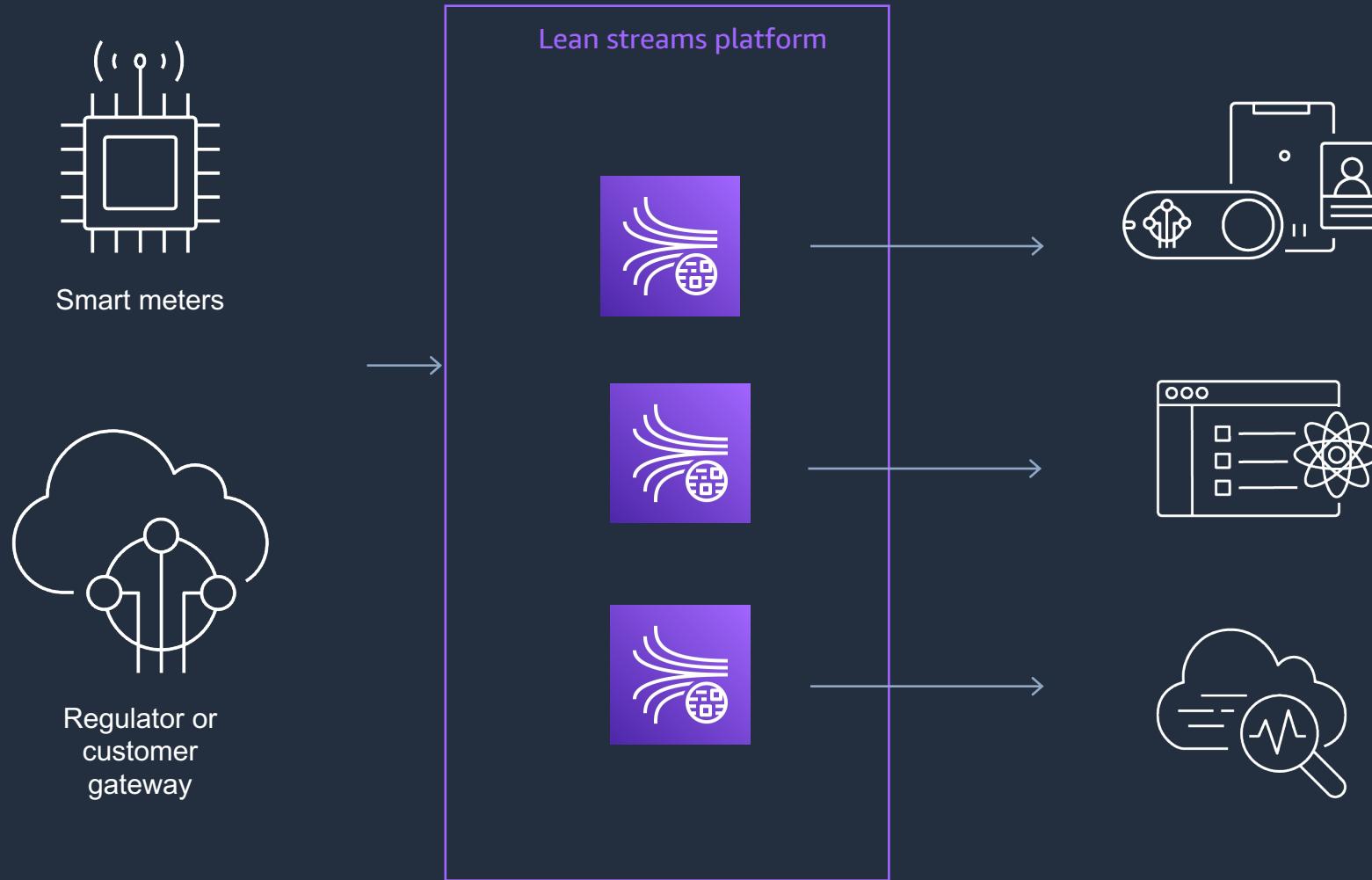
Both architectures are not designed  
to solve challenges



# Working backwards



# Working backwards

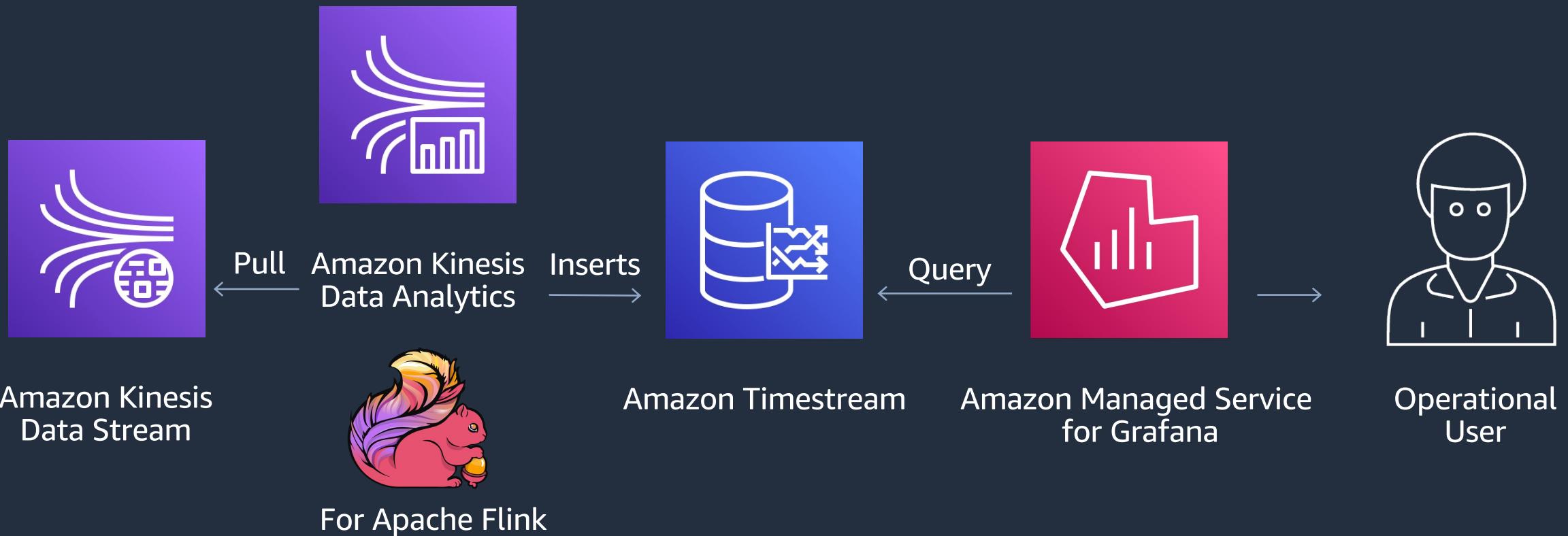


# Speed consumption



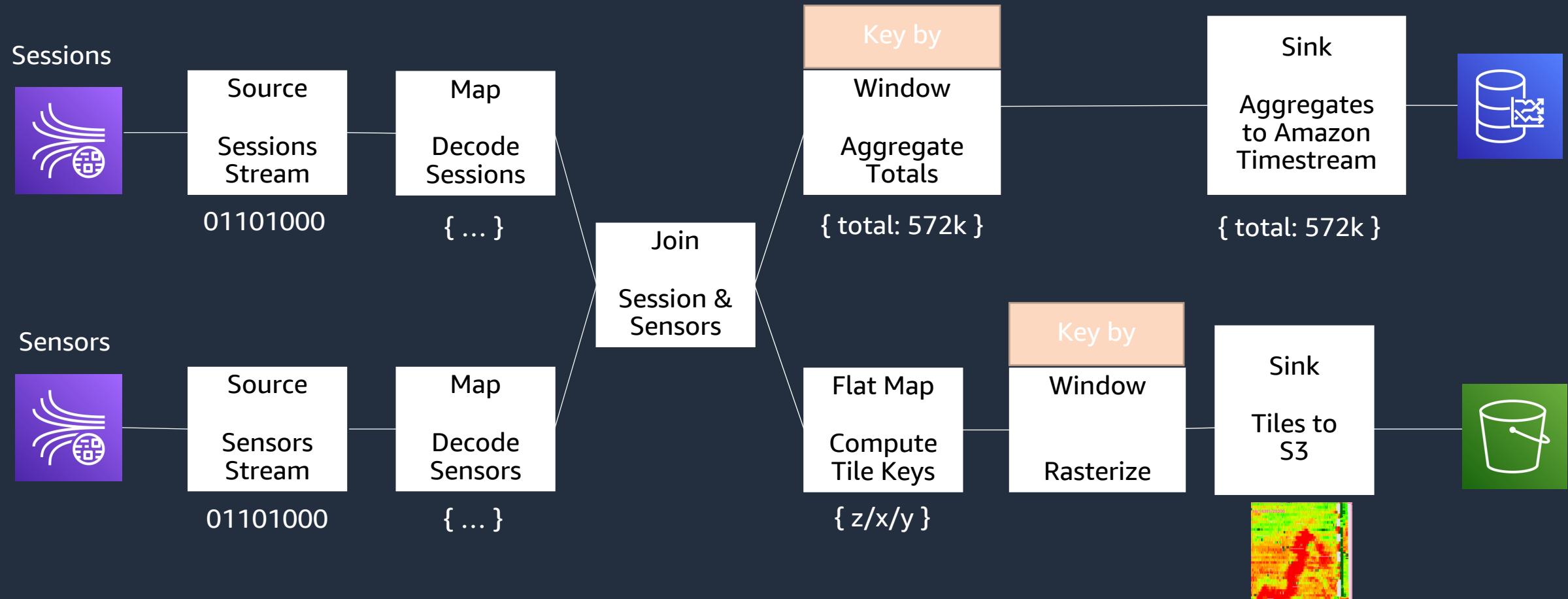
© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

# Near-Real Time Operational Insights Architecture



# Streaming Analytics

APACHE FLINK - EXAMPLE



# Streaming Analytics

APACHE FLINK – DECLARATIVE STREAMS

```
createKinesisSource(env, parameter)
    .map(JsonToTimestreamPayloadFn()).name("MaptoTimestreamPayload")
    .process(OffsetFutureTimestreamPoints()).name("UpdateFutureOffsetedTimestreamPoints")
    .addSink(TimestreamSink(region, databaseName, tableName, batchSize))
    .name("TimestreamSink<$databaseName, $tableName>")
```

# Amazon Timestream Architecture

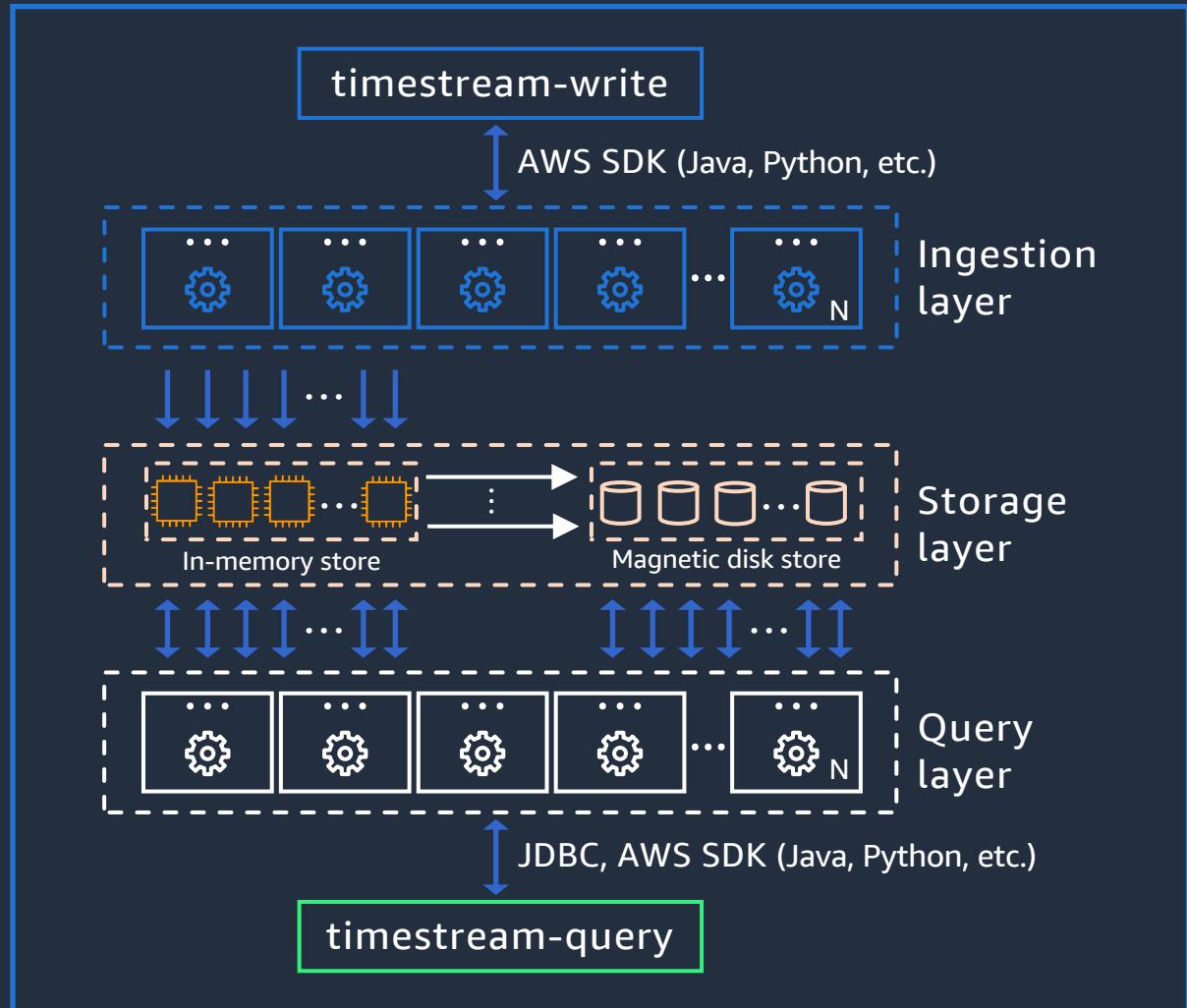
PURPOSE-BUILT FOR TIME SERIES WORKLOADS

## Decoupled architecture

- Highly available
- Independently scalable ingestion, storage, and query processing

## High throughput auto scaling ingestion

- Durable: data is replicated across multiple Availability Zones
- Automatic data deduplication handling
- No need to provision or configure resources



# Amazon Timestream Architecture

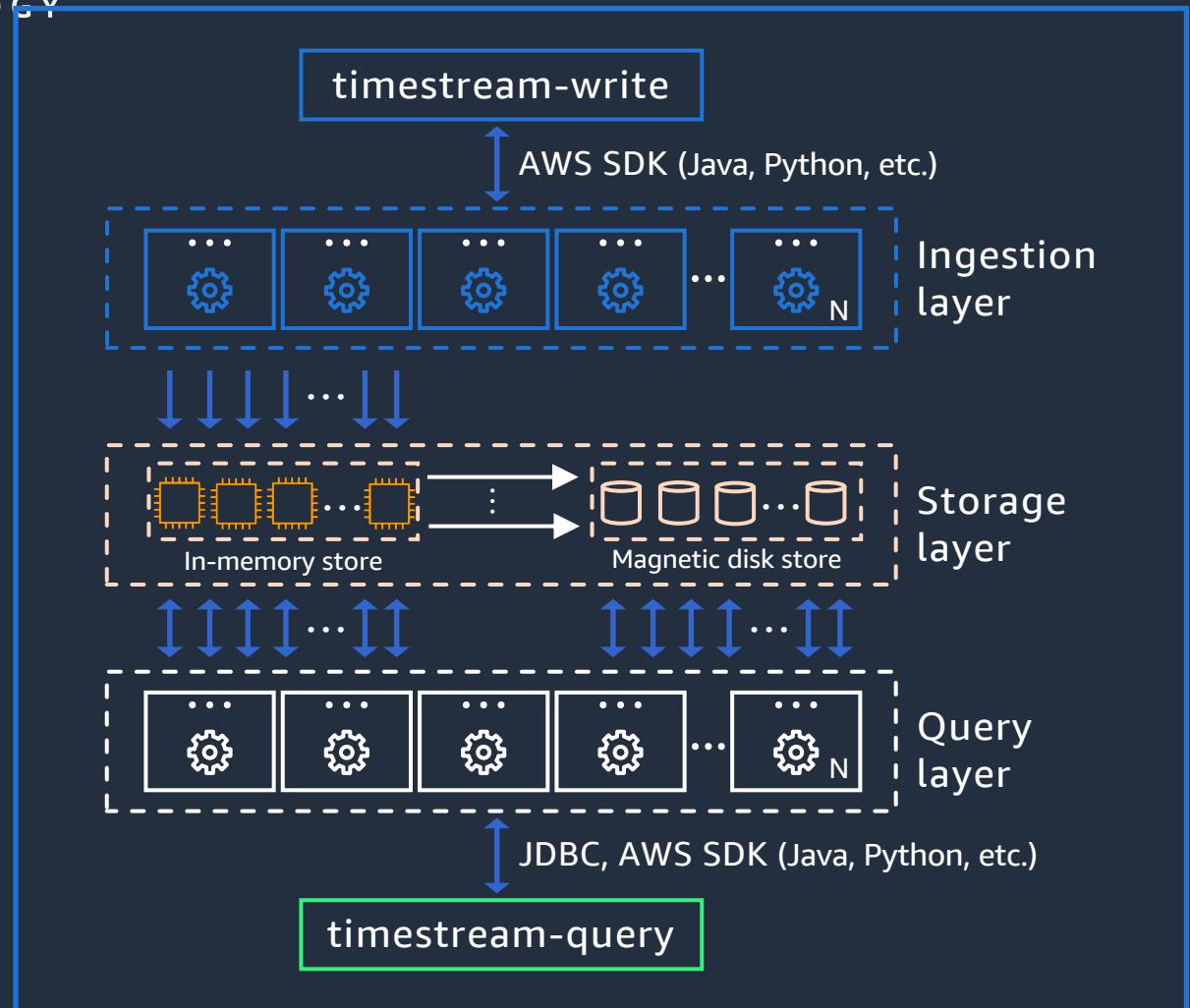
SERVERLESS, CLOUD-NATIVE AWS BUILT TECHNOLOGY

## Multiple tiers of storage

- Scalable to petabytes and beyond
- **Memory store**: fast point-in-time queries, high-throughput durable ingestion
- **Magnetic store**: high-performance analytical queries and low-cost long-term storage

## Scalable SQL query engine

- **Adaptive query engine** is capable of querying data across multiple data tiers
- **No indexes** to configure and **no provisioning** required



# Operational Views

AMAZON MANAGED SERVICE FOR GRAFANA

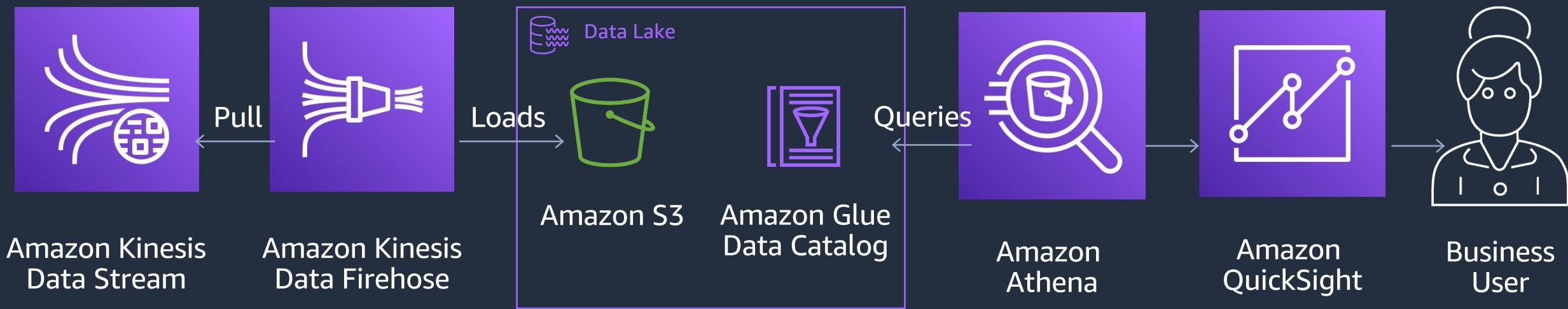


# Big data insights



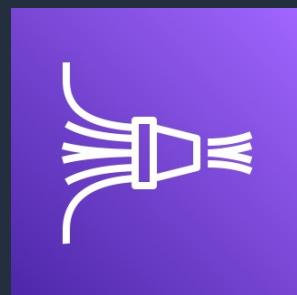
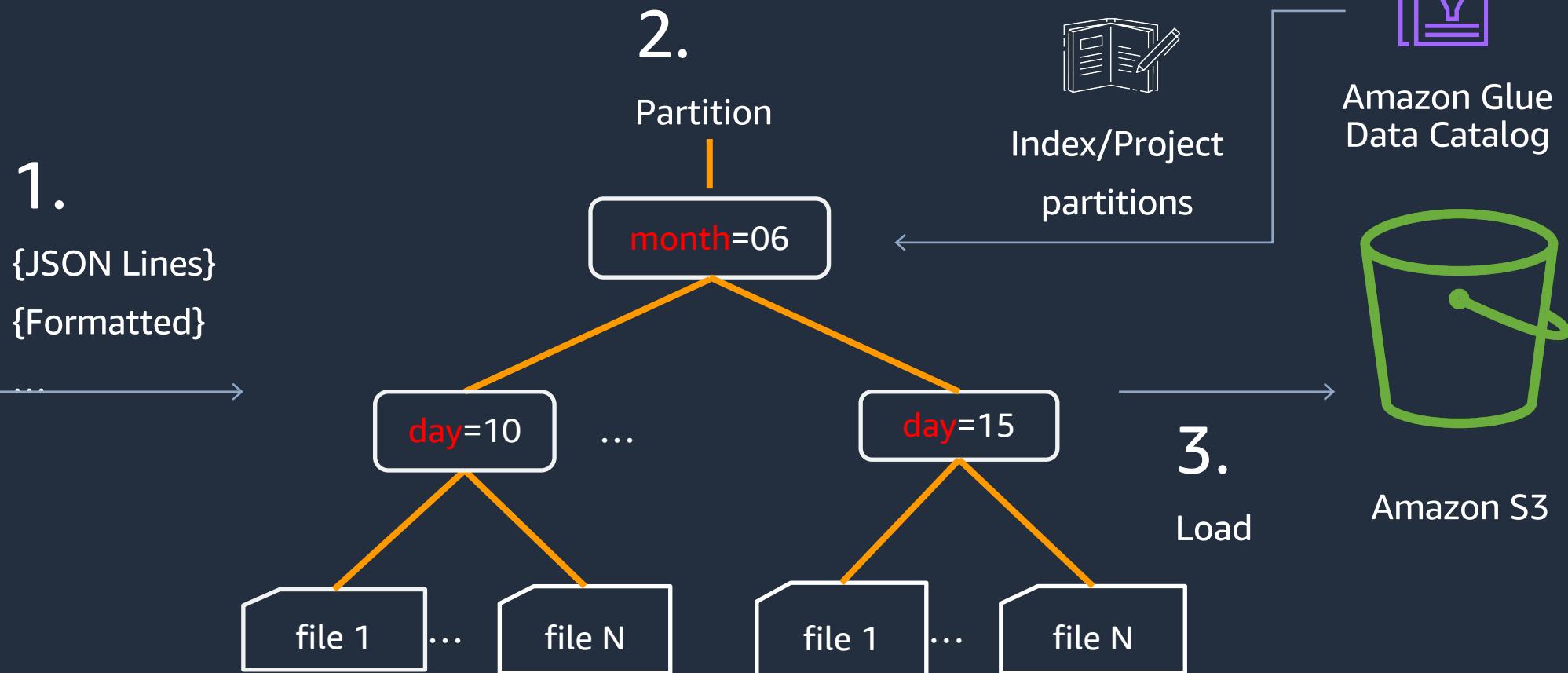
© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

# Strategic Big Data Insights Architecture

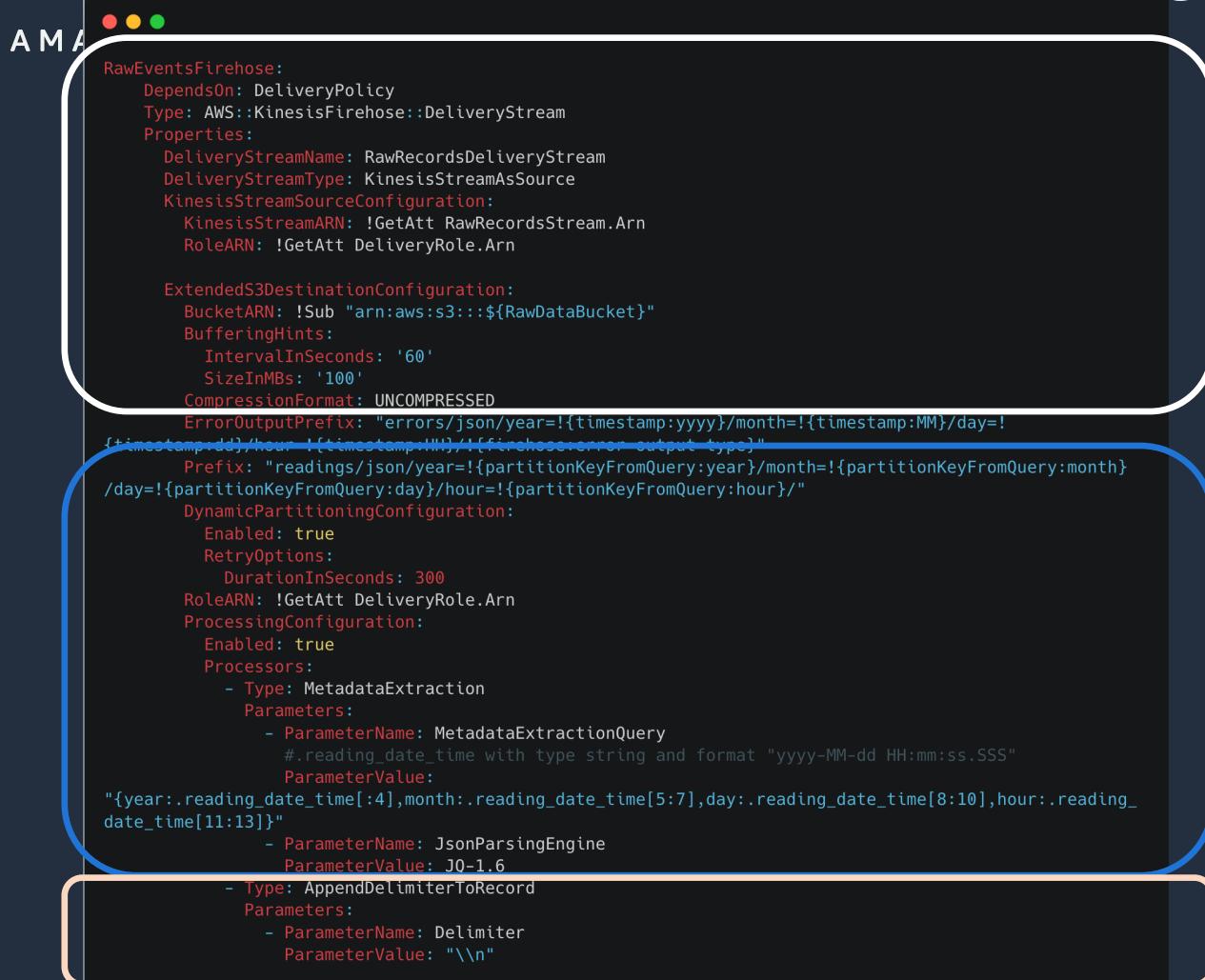


# Data Encoding and Partitioning

AMAZON KINESIS DATA FIREHOSE

Amazon Kinesis  
Data Firehose

# Data Encoding and Partitioning



Move data from Stream to Raw S3 Bucket

Late arriving data compensation

JSON lines

# Discover Dynamic Partitions

## GLUE DATA CATALOG PARTITION PROJECTION

Prefix: "readings/json/year=!{partitionKeyFromQuery:year}/month=!{partitionKeyFromQuery:month}/day=!{partitionKeyFromQuery:day}/hour=!{partitionKeyFromQuery:hour}"/"

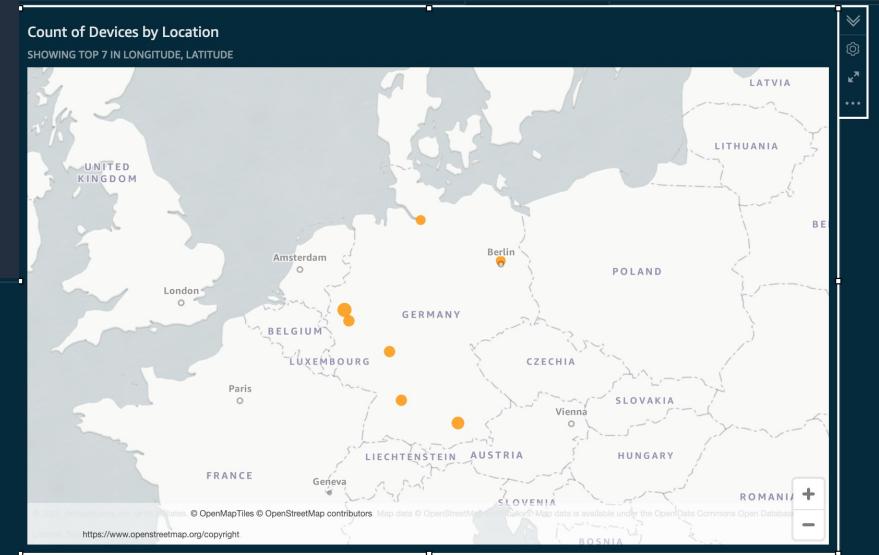
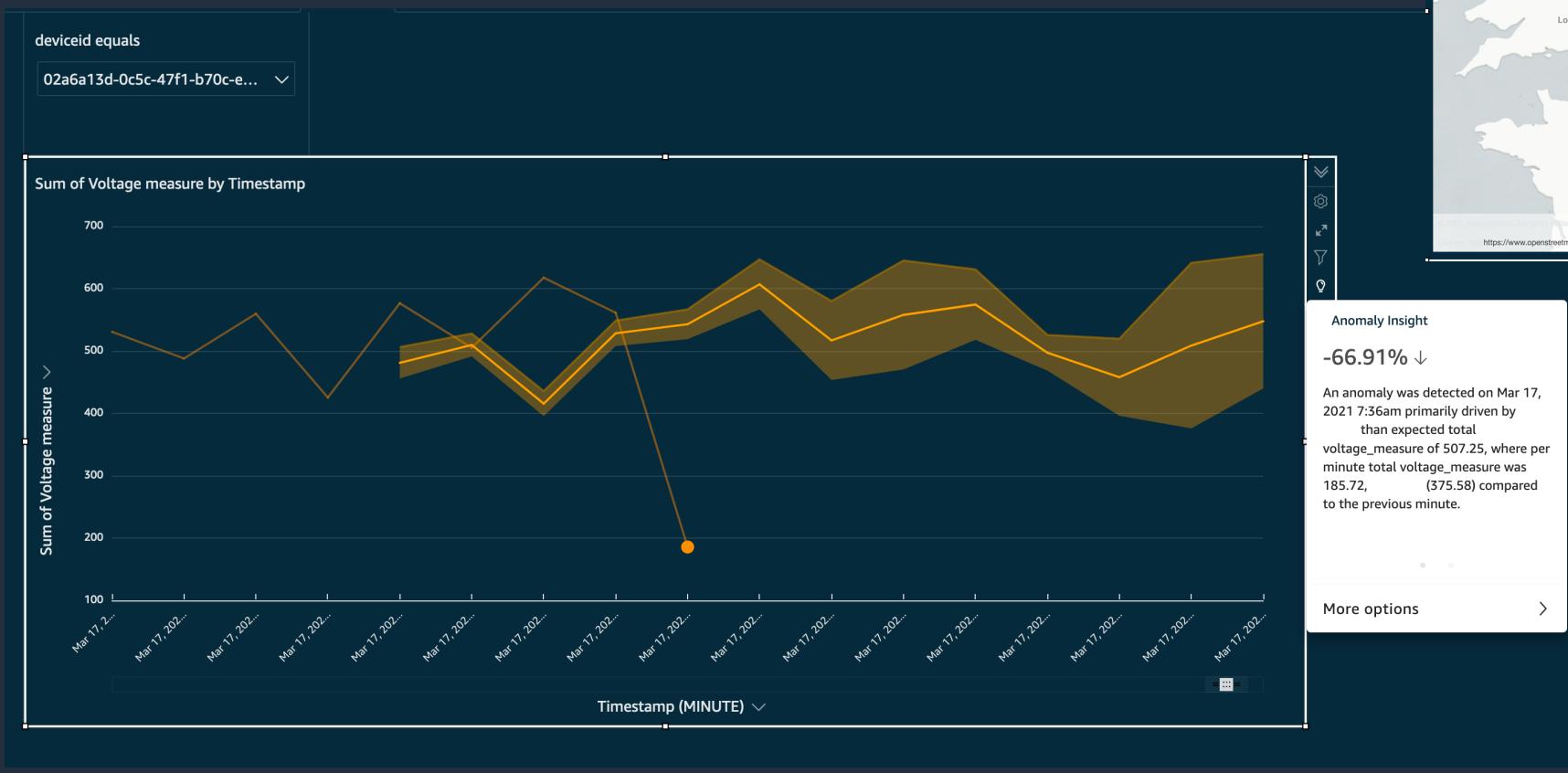
bucket\_name > measures/ > json/ > year=2023/ > month=01/ > day=02/ > hour=01/

```
CREATE EXTERNAL TABLE raw_device_data (
  ...
)
...
PARTITIONED BY (
  year INT,
  month INT,
  day INT,
  hour INT,
)
LOCATION "s3://DOC-EXAMPLE-BUCKET/prefix/"
TBLPROPERTIES (
  "projection.enabled" = "true",
  "projection.year.type" = "integer",
  "projection.year.range" = "2011,2099",
  "projection.year.digits" = "4",
  "projection.month.type" = "integer",
  "projection.month.range" = "1,12",
  "projection.month.digits" = "2",
  "projection.day.type" = "integer",
  "projection.day.range" = "1,31",
  "projection.day.digits" = "2",
  "projection.hour.type" = "integer",
  "projection.hour.range" = "0,23",
  "projection.hour.digits" = "2",
  "storage.location.template" = "s3://DOC-EXAMPLE-BUCKET/prefix/year=${year}/month=${month}/day=${day}
/hour=${hour}/"
```

"projection.enabled" = "true",  
"projection.year.type" = "integer",  
"projection.year.range" = "2011,2099",  
"projection.year.digits" = "4",

# Insights Views

AMAZON QUICKSIGHT

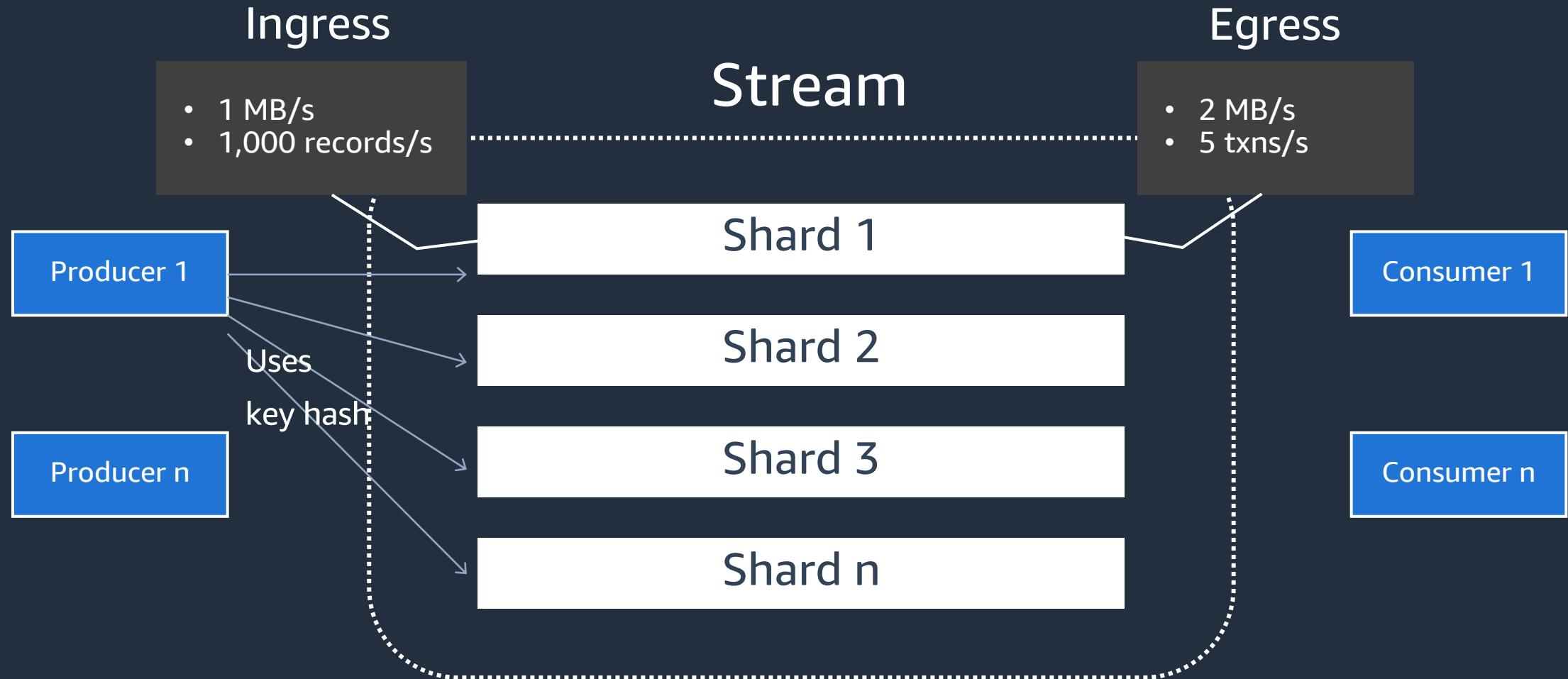


# Data stream



© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

# Kinesis Data Streams: Granular Provisioning and Scaling



# Scaling mechanism: Default and max

New on-demand stream

Writes: 4 MB/s and  
4,000 records/s

Reads: 8 MB/s

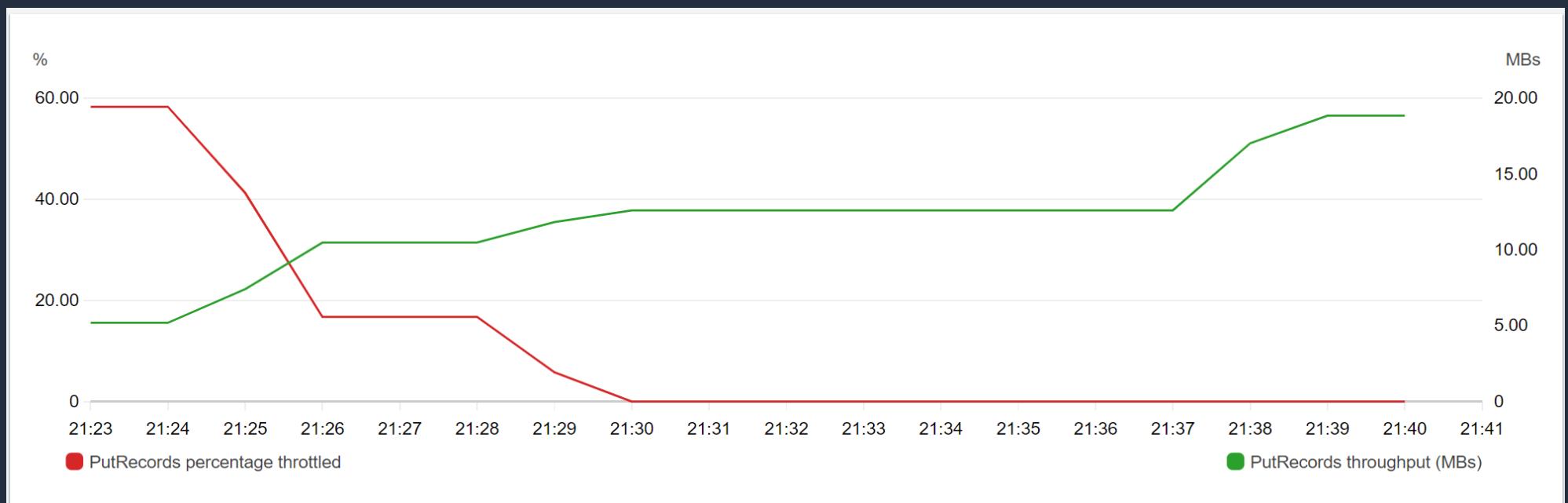
Maximum

Writes: 200 MB/s  
and 200,000  
records/s

Reads: 400 MB/s

# Scaling mechanism in action

- At time stamp 21:23, we change the stream capacity from provisioned to on-demand
- The stream immediately scales to 10 shards and throttling is reduced
- Now the traffic continues to increase to 18 MB/s, so KDS scales to 40 shards by 21:30 and there are no throttles
- You can also see that the throttles start to reduce within minutes after moving the stream into on-demand and are fully eliminated even as traffic increases



# Ingestion

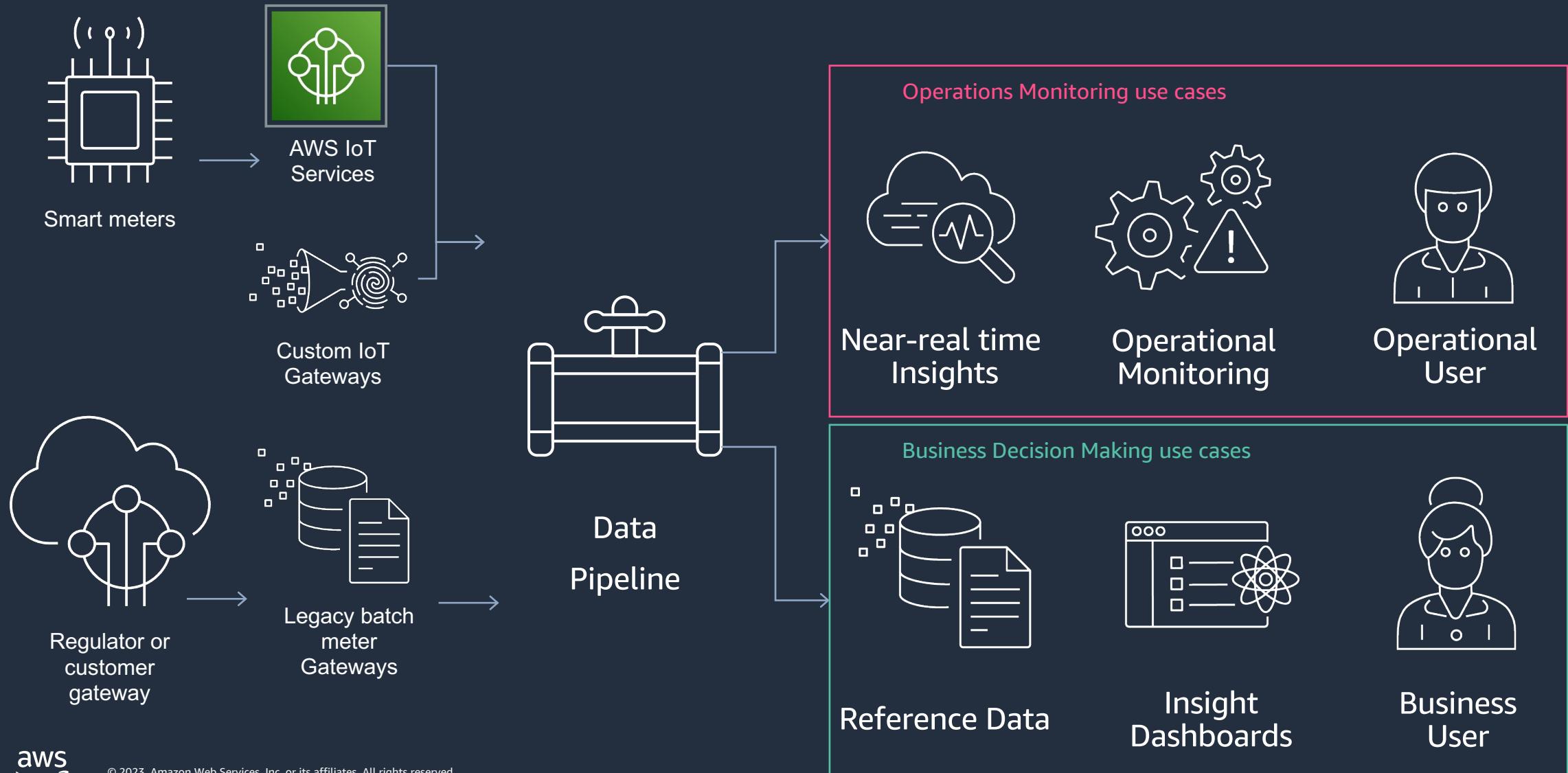


© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

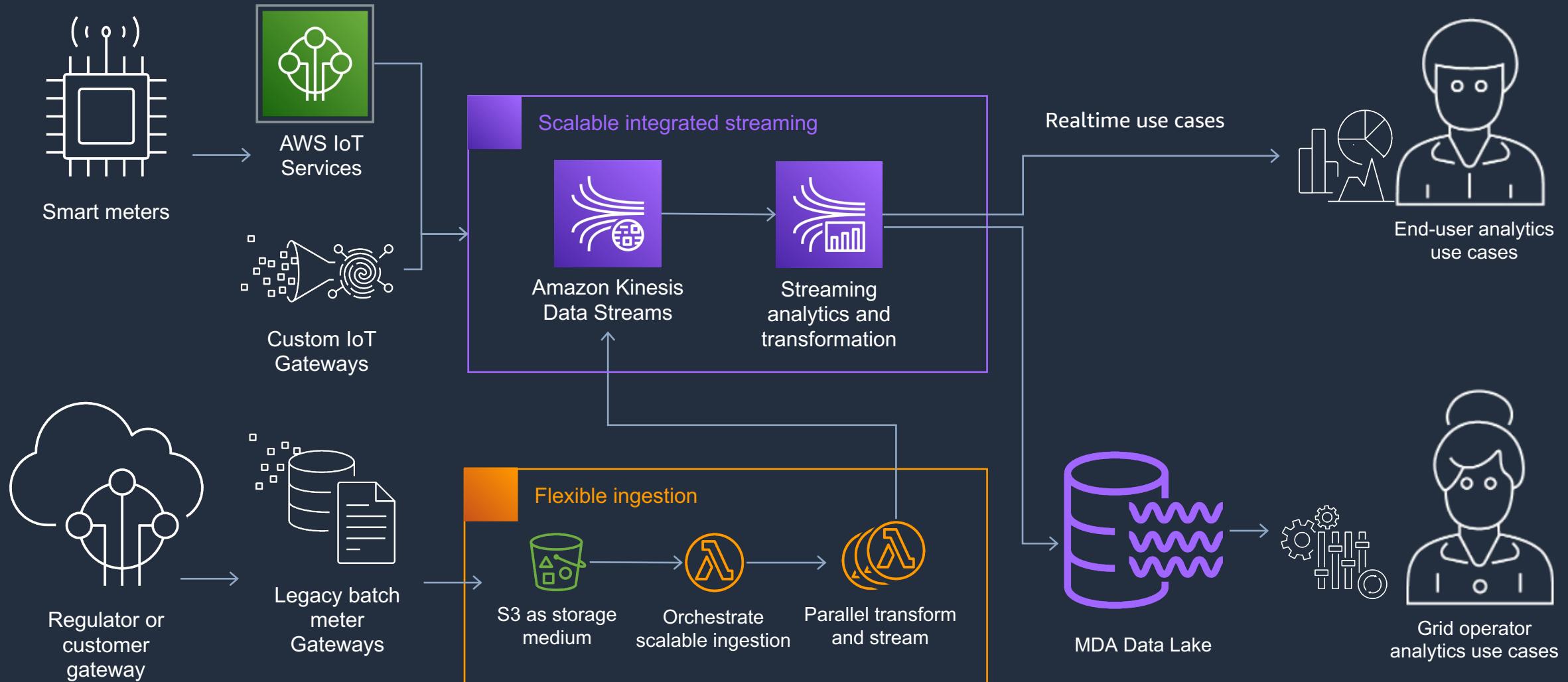


# Your stream should not scale to all kinds of ingestion

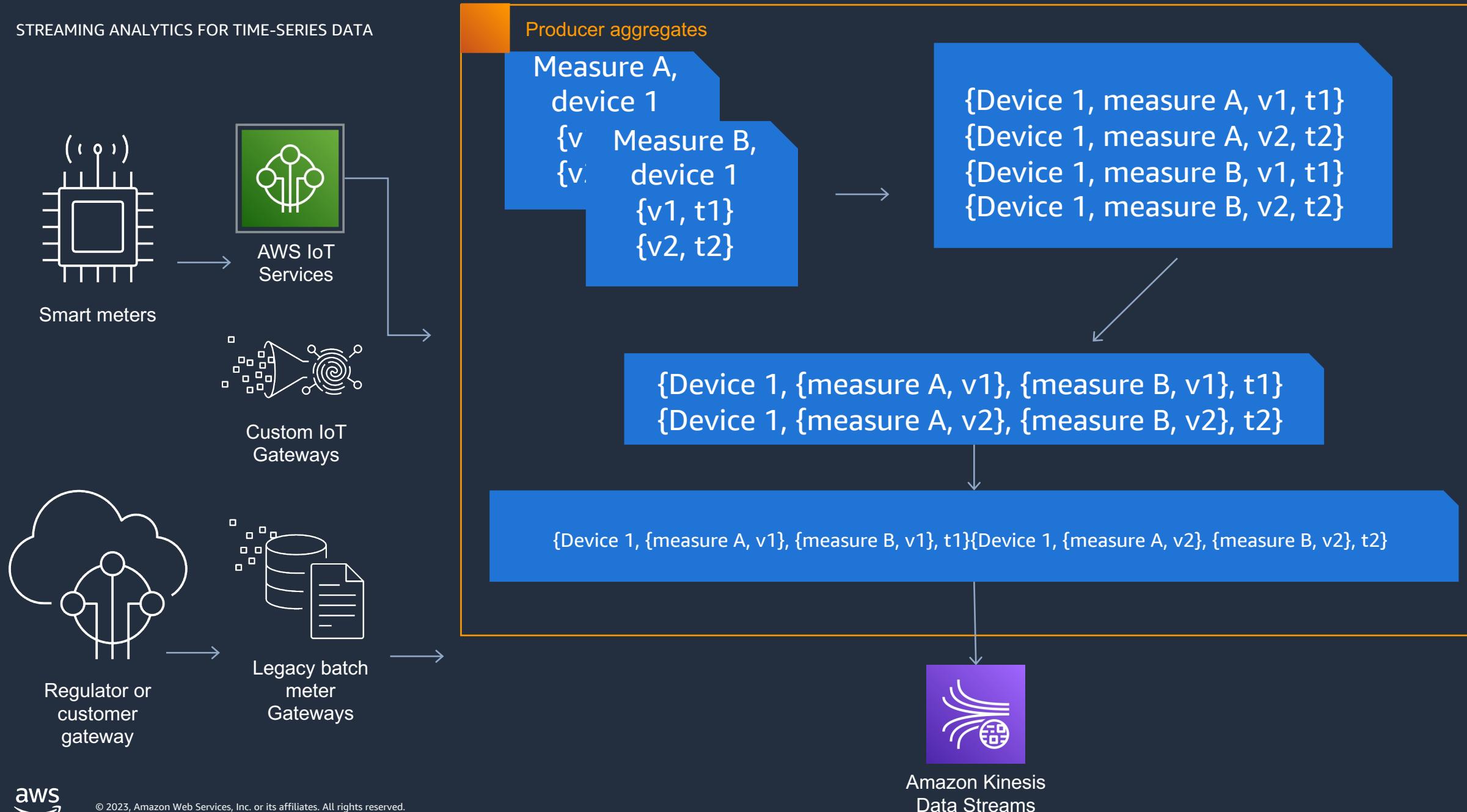
## STREAMING ANALYTICS FOR TIME-SERIES DATA



## STREAMING ANALYTICS FOR TIME-SERIES DATA



## STREAMING ANALYTICS FOR TIME-SERIES DATA



```
agg = SimpleJsonRecordAggregator()  
agg.on_record_complete(callback=consume_staging_records_intake, executor= map_async)  
  
response = s3_select(bucket, prefix, "SELECT * FROM S3Object", start_range, end_range)  
"""  
    select_object_content() response is an event stream that can be looped to concatenate the overall result set  
    Hence, we are joining the results of the stream in a string before converting it to a tuple of dict  
"""  
for index, payload in enumerate(response["Payload"]):  
    if records := payload.get("Records"):  
        staging_records = (records["Payload"].decode("utf-8")).split("\n")  
        for some_staging_record in staging_records:  
            agg.add_record(some_staging_record["meter_id"], some_staging_record)  
agg.flush()
```

Reactively producer to stream

Stream from S3

Add records to aggregator

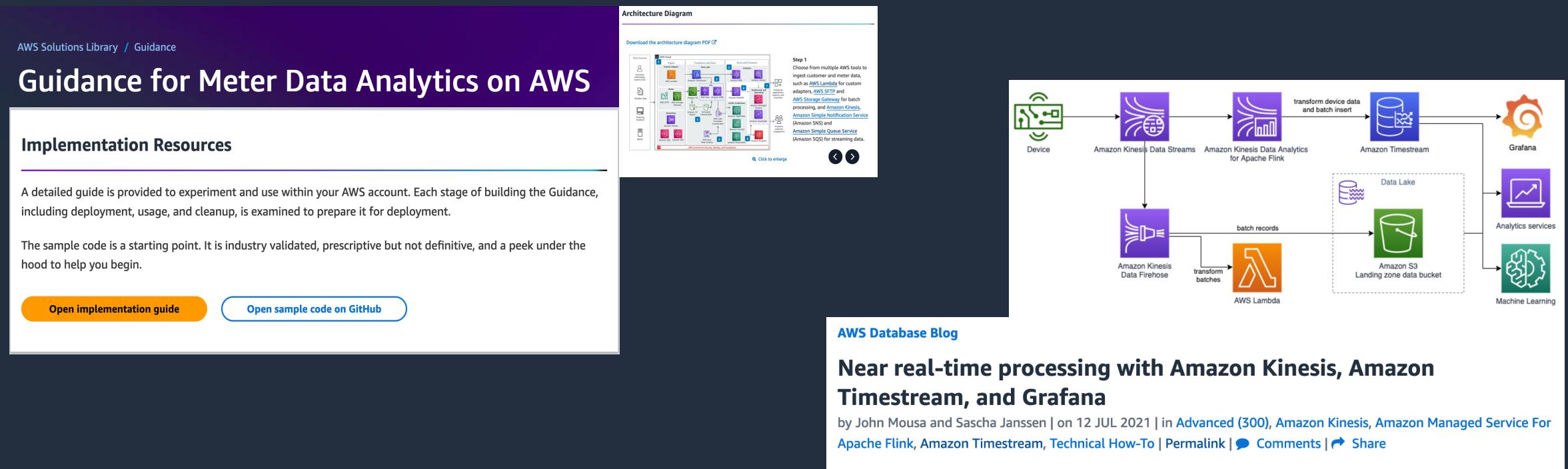
# Summary



© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

# Getting started

<https://aws.amazon.com/solutions/guidance/meter-data-analytics-on-aws/>



<https://github.com/aws-samples/amazon-kinesis-timestream-grafana>



<https://www.pulse.aws/survey/FL1FAATO>

# Thank you!

John Mousa



Direct link to slides right after 5-click survey feedback

















“

”





# Thank you!