



Streaming data pipelines for real-time analytics

– Are you ready?

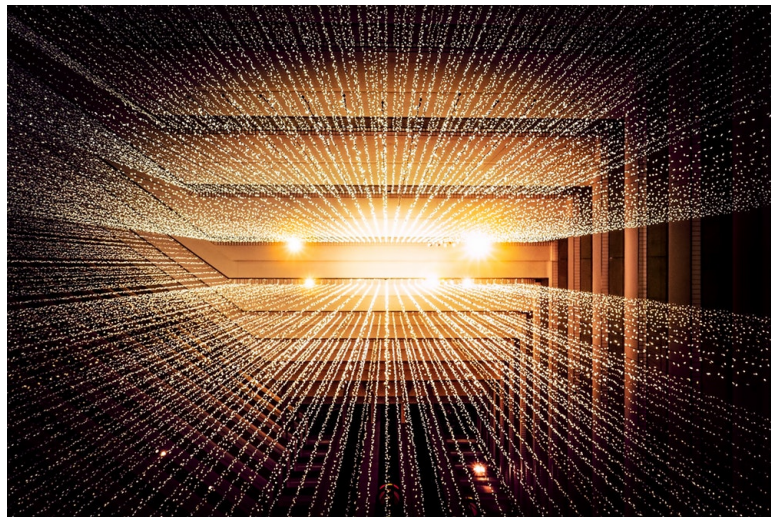
Low latency and low cost at any scale

Sai Maddali, Senior Product Manager, Amazon Kinesis



Table of contents

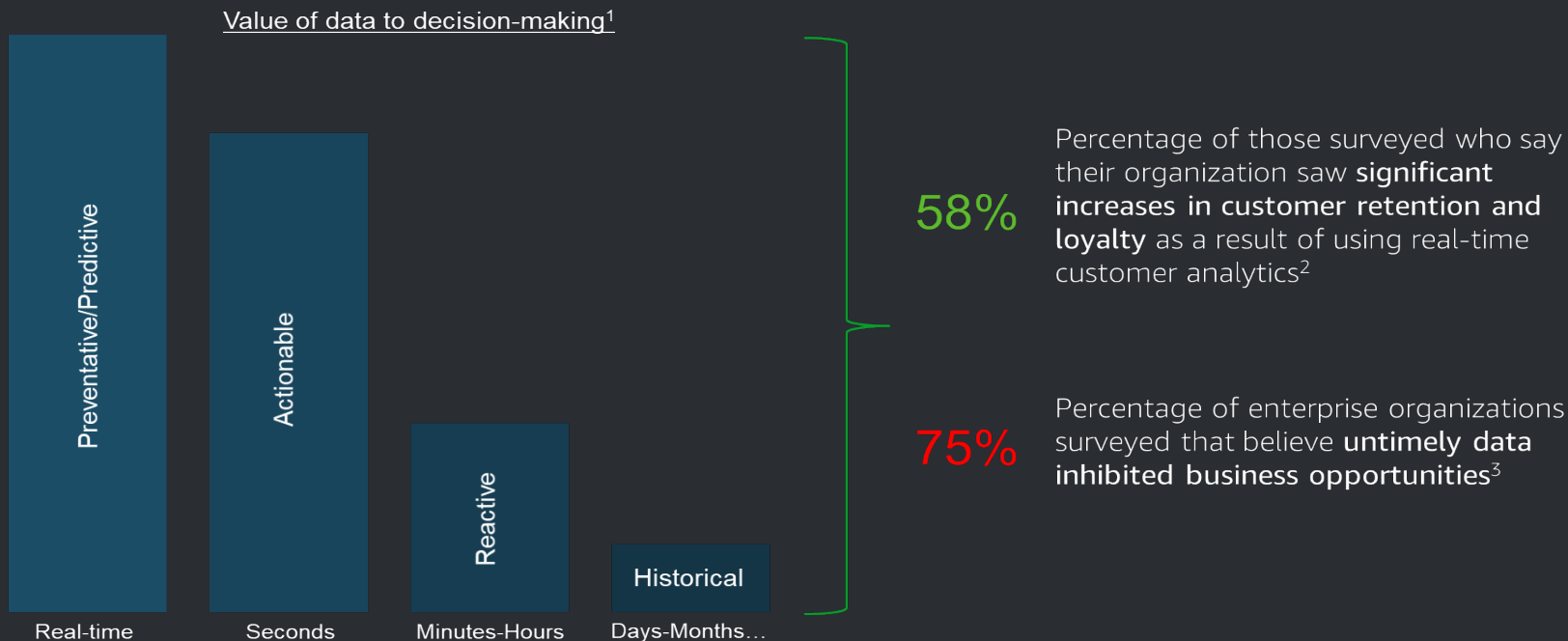
- Why data streaming?
- 5 Kinesis super powers
- Use cases to get started
- Questions



Why data streaming?

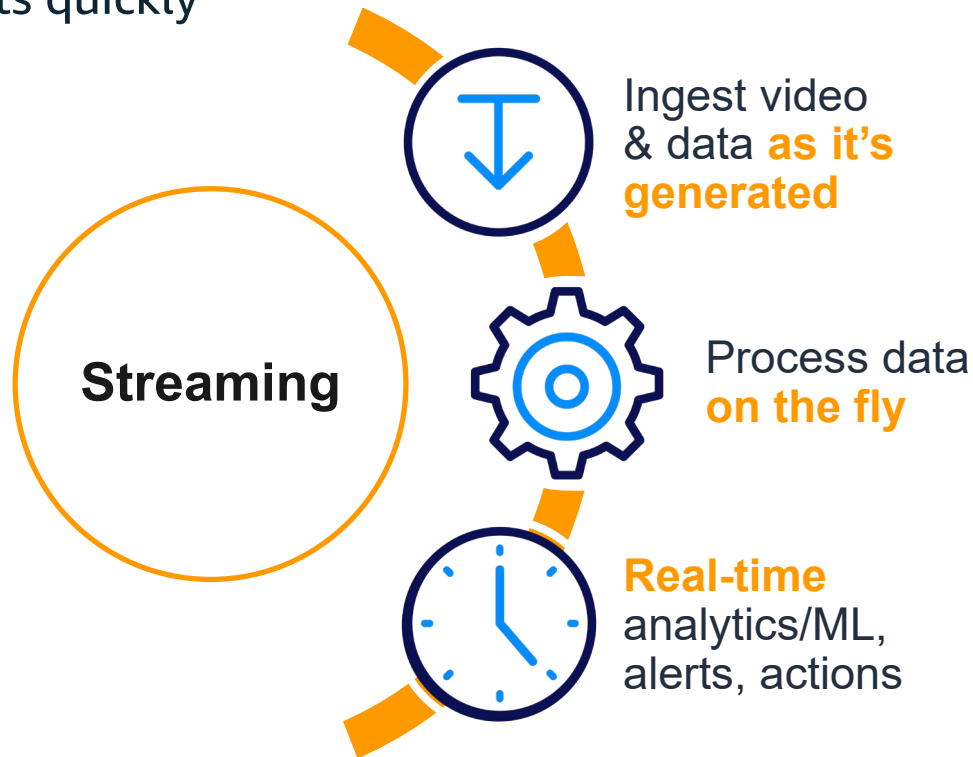
Timely decisions require new data in minutes

Data have a short shelf life of actionability¹. AWS lets you act on that data **as fast as the market dictates**.



Stream new data in seconds

Get actionable insights quickly



Epic Games continually improves Fortnite for 250+ million players globally

Challenge:

They needed a way to process and analyze over **100PB** of data (**125M events/min**) ingested from game clients and game servers to understand and adapt to player engagement.

Solution:

Epic Games turned to AWS for an Amazon S3 data lake in combination with Amazon EMR, Amazon EC2, and Amazon Kinesis.

Result:

The data provides a constant feedback loop for designers, and an up to the minute analysis of gamer satisfaction to drive gamer engagement.



Amazon S3



Amazon EMR



Amazon Kinesis



Amazon EC2



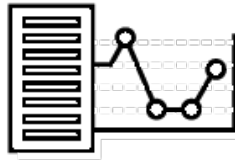
Most common uses of streaming



Security
Monitoring



Industrial
Automation



Log
Analytics



Data
Lakes



Microservices
communication

Enabling real-time analytics

Data streaming technology enables a customer to ingest, process, and analyze high volumes of high-velocity data from a variety of sources in real time



Source

Devices and or applications that produce real-time data at high velocity

Stream ingestion

Data from tens of thousands of data sources can be written to a single stream

Stream storage

Data is stored in the order it was received for a set duration of time, and it can be replayed indefinitely during this time

Stream processing

Records are read in the order they are produced, enabling real-time analytics or streaming ETL

Destination

Data lake (most common)
Analytics services
Database (least common)

Amazon Kinesis: Real-time streaming on AWS

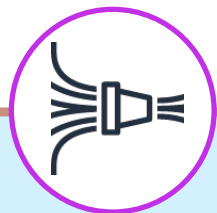
Easily collect, process, and analyze data streams in real time

**Amazon Kinesis
Data Streams**



Collect and store data streams for analytics

**Amazon Kinesis
Data Firehose**



Load data streams into data stores

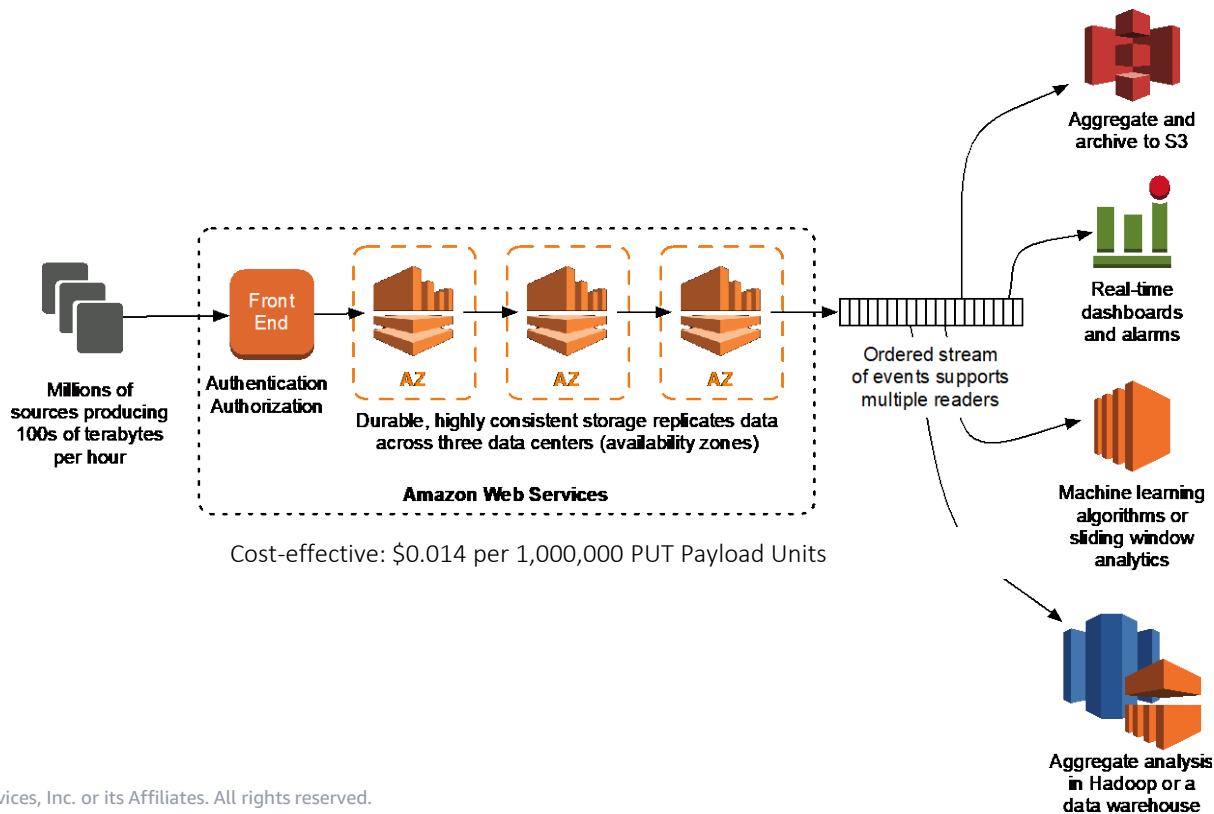
**Amazon Kinesis
Data Analytics**



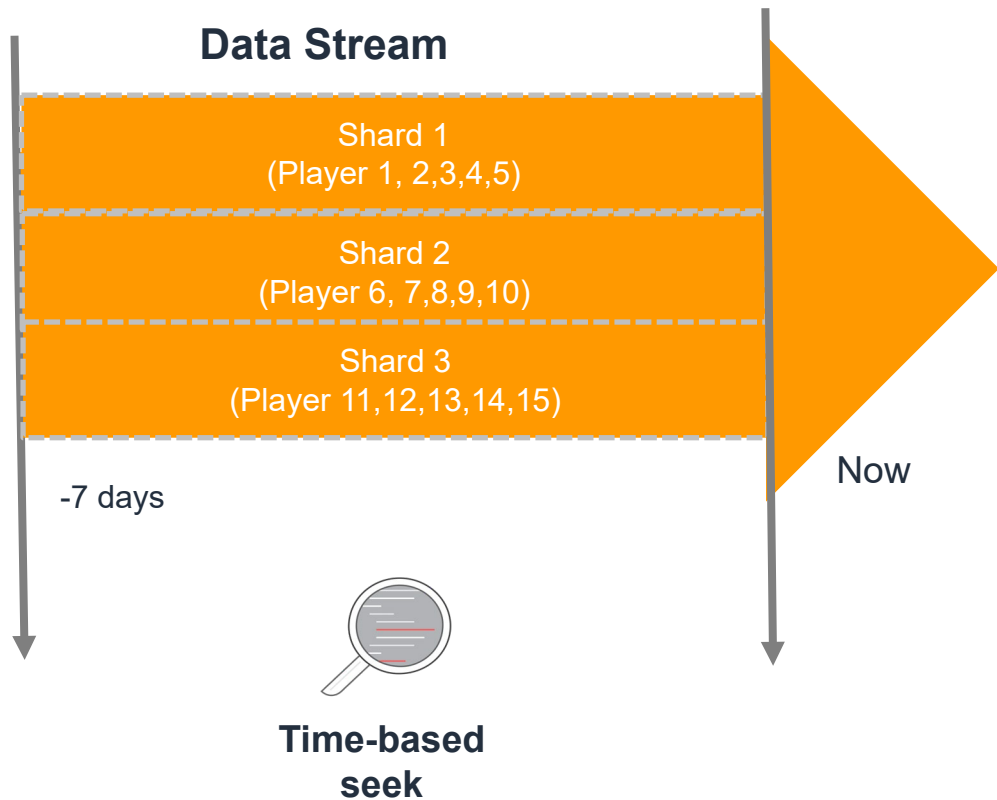
Analyze data streams with SQL or Java (Apache Flink)

Amazon Kinesis Data Streams – How it works

Fully managed service for real-time processing of streaming data



Managed ability to capture and store data



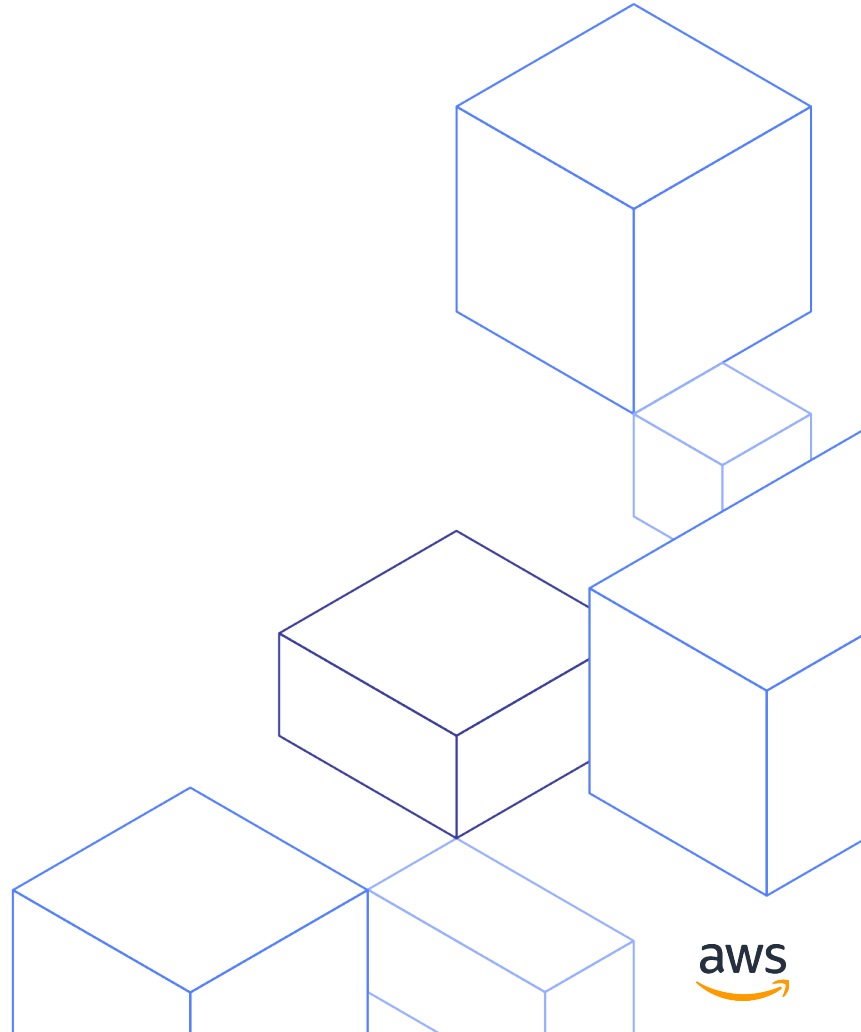
- Data streams are made of **Shards**
- **Shard** – Unit of throughput and parallelism
- **Partition key** - Business key that enables automatic data mapping into shards.
- **Iterator** - Able to seek at any point in the stream to read data

5 super powers of Amazon Kinesis Data Streams



1. Easy to get started
2. Easy to operate
3. Massive scale
4. Low-latency
5. Low cost

Easy to get started



Get started in minutes with a few clicks

aws Services Resource Groups

Data stream name
MyMassiveStream
Acceptable characters are uppercase and lowercase letters, numbers, underscores, hyphens and periods.

Data stream capacity [Info](#) [Request limit increase](#)

Data records are stored in Kinesis Data Stream. A shard is a uniquely identified sequence of data records in a stream.

▶ Shard estimator

Number of open shards
Each shard ingests up to 1 MiB/second and 1000 records/second and emits up to 2 MiB/second.

100

Minimum: 1, Maximum: 489, Account limit: 500.

Total data stream capacity
Total data stream capacity is calculated based on the number of shards entered above.

Write
100 MiB/second, 100000 Data records/second

Read
200 MiB/second

Cancel **Create data stream**

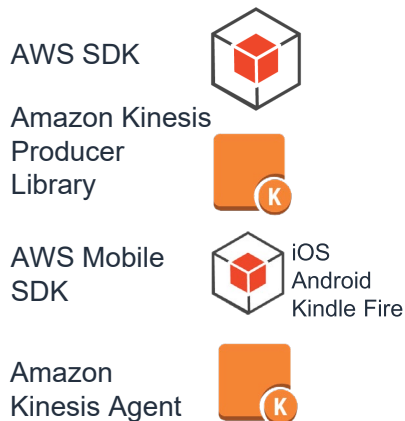
Feedback English (US)

- Guaranteed throughput, making it easy to size the workload
- Create a data stream with the necessary throughput in minutes

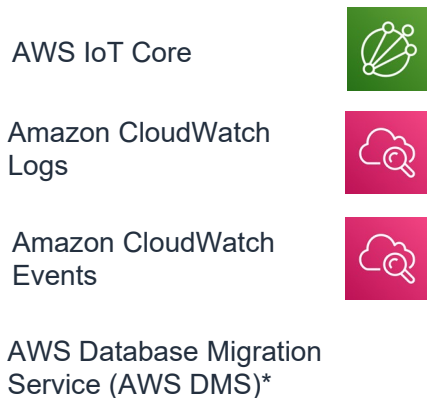
Integrate with existing systems to ingest data

Data from tens of thousands of data sources can be written to a single stream

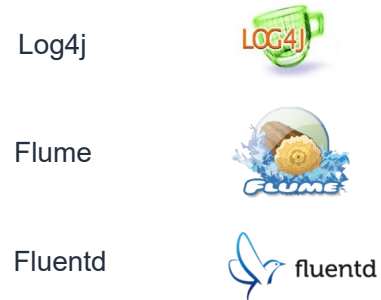
AWS toolkits/libraries



AWS service integrations

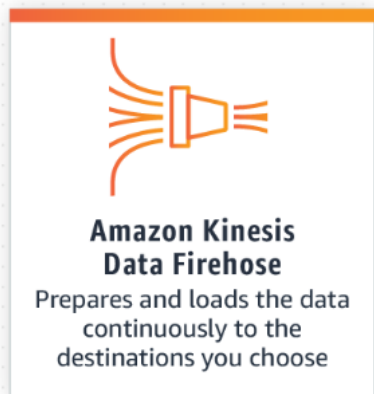


Third-party offerings



*AWS DMS includes eight on-premises databases, one Azure database, five Amazon RDS/Amazon Aurora database types, and Amazon Simple Storage Service (Amazon S3)

Deliver data to different destinations in few clicks



Amazon Kinesis Data Firehose

Deliver data to destinations such as Amazon S3, Amazon Redshift, Amazon Elasticsearch Service, and generic **HTTP endpoints** so you can use existing analytics tools

Write data processing applications quickly



Use AWS Lambda to quickly process streaming data



Get actionable insights from streaming data in real time

```
Add SQL from templates Download SQL
```

```
1 CREATE STREAM sliding_window (device_parameter VARCHAR(16), sum_device_value INTEGER, record_count_in_window INTEGER);
2
3 CREATE PUMP sliding_pump AS INSERT INTO sliding_window
4 SELECT STREAM device_parameter, max(device_value) OVER w1, count(*) OVER w2 as record_count_in_window
5 FROM source_sql_stream_001
6 WINDOW w1 AS (PARTITION BY device_parameter RANGE INTERVAL '1' MINUTE PRECEDING);
7
8 CREATE STREAM max_window (device_parameter VARCHAR(16), max_count INTEGER);
9
10 CREATE PUMP max_pump AS INSERT INTO max_window
11 SELECT STREAM device_parameter, max(record_count_in_window) as max_count
12 FROM sliding_window
13 GROUP BY device_parameter, STEP(sliding_window.routine BY INTERVAL '5' SECOND);
14
```

Exit (done editing) Save and run SQL



Kinesis Data Analytics for Java for sophisticated applications

Uses Apache Flink, a framework and distributed engine for stateful processing of data streams



Simple programming

Easy-to-use and flexible APIs make building apps fast



High performance

In-memory computing provides low latency & high throughput



Stateful processing

Durable application state saves



Strong data integrity

Exactly-once processing and consistent state

Native AWS Integrations



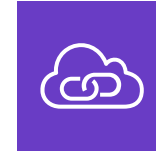
AWS Identity and Access Management



AWS Key Management Service



Amazon VPC



AWS PrivateLink



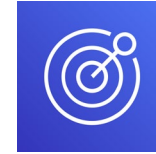
Amazon API Gateway



Amazon CloudWatch



Amazon EventBridge



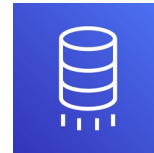
Amazon Pinpoint



Amazon Redshift



Amazon Quantum Ledger Database

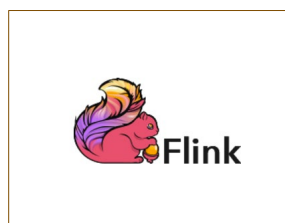
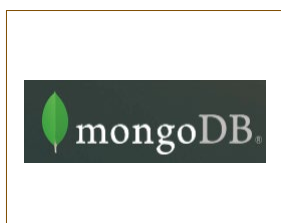
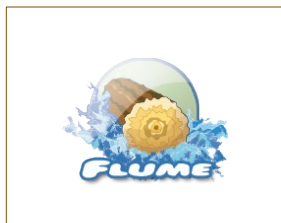
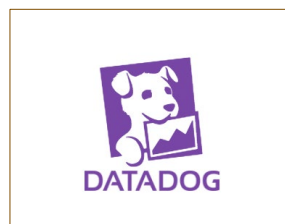


AWS Database Migration Service

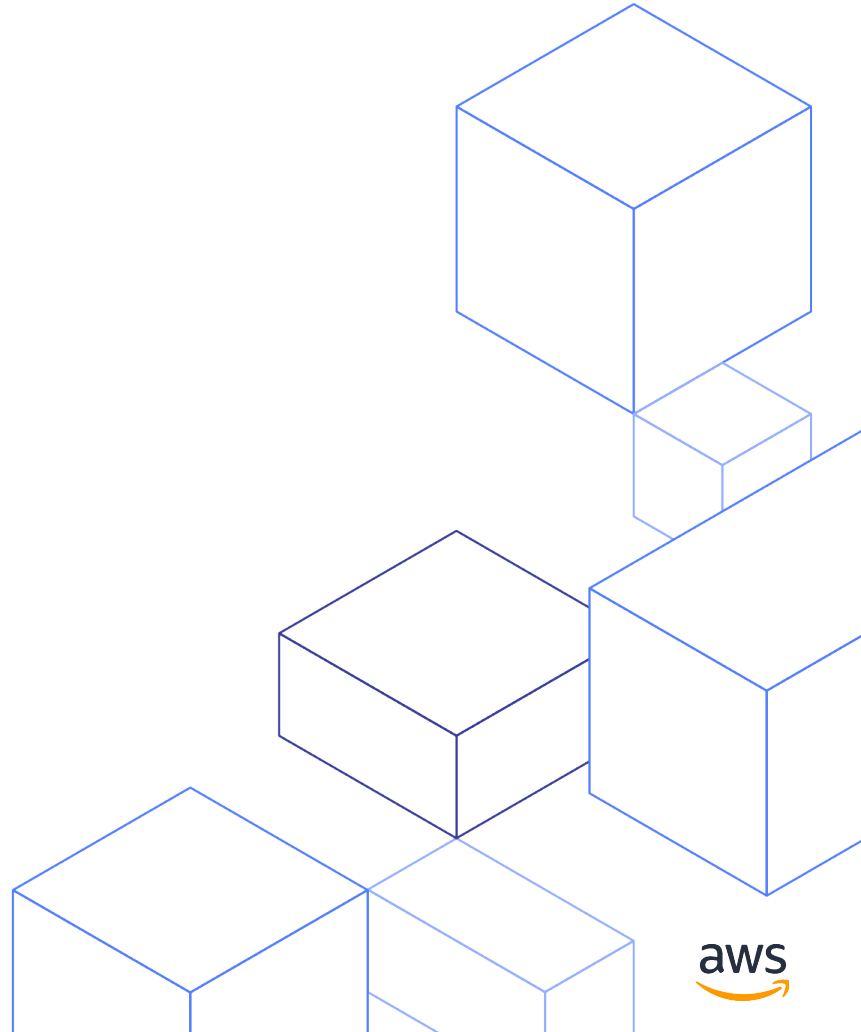


Amazon Elasticsearch Service

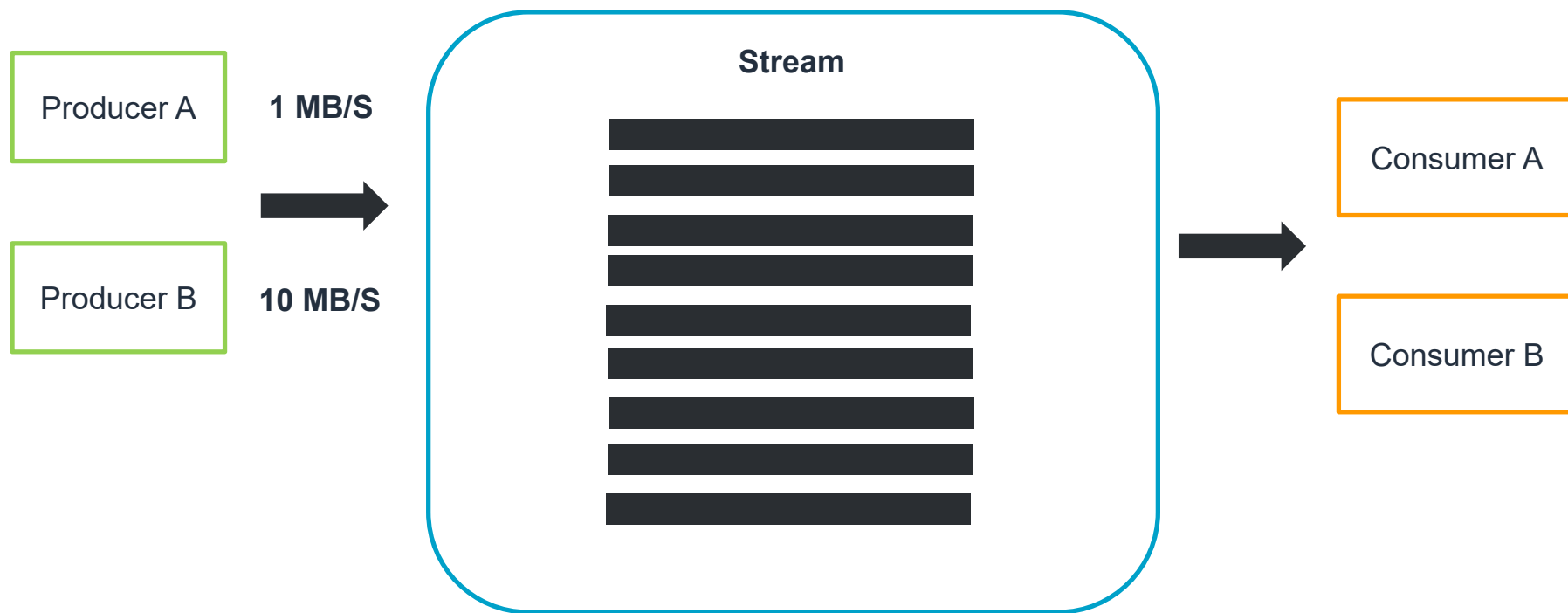
3rd Party Connectors



Easy to operate

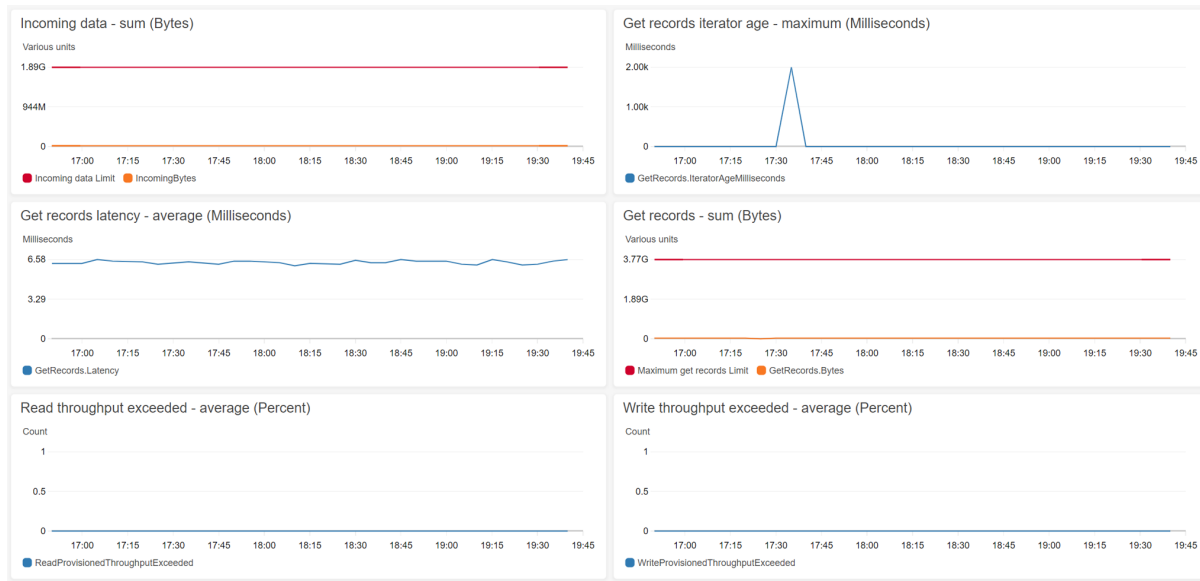


Seamless and non-disruptive scaling

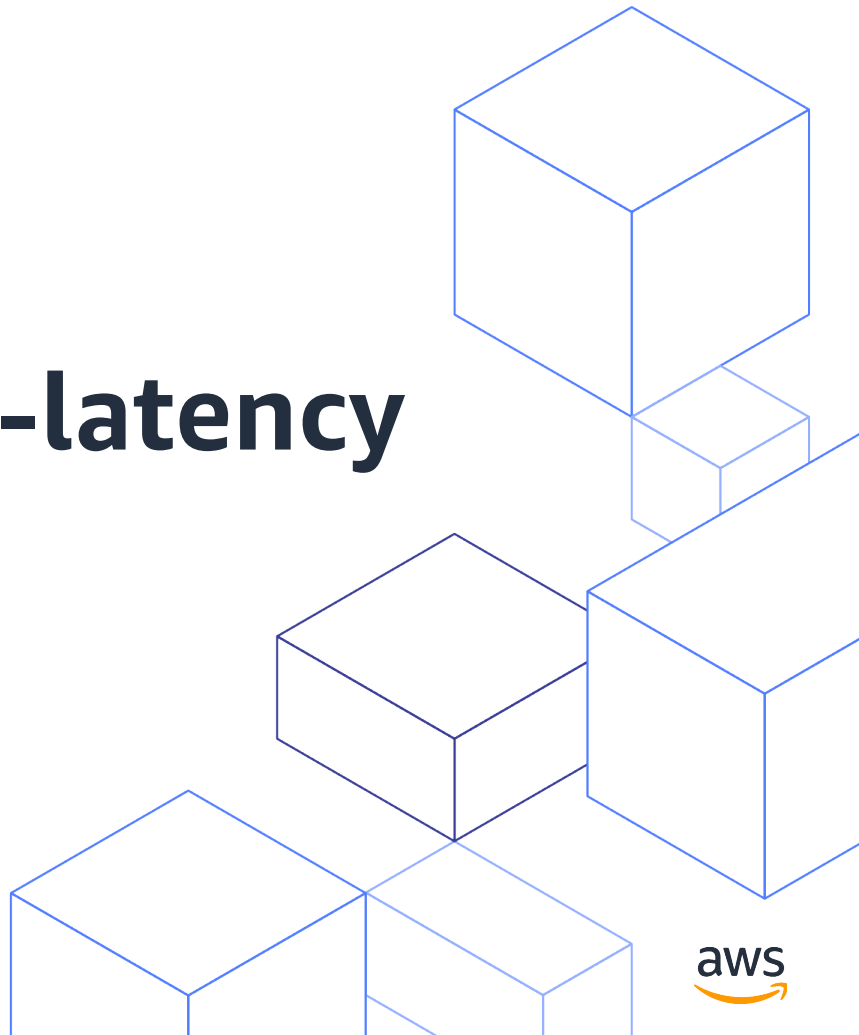


Low operations overhead to manage large streams

- Quickly identify and remediate issues using integration with Amazon CloudWatch
- Automate capacity management using CloudWatch and AWS Auto-scaling

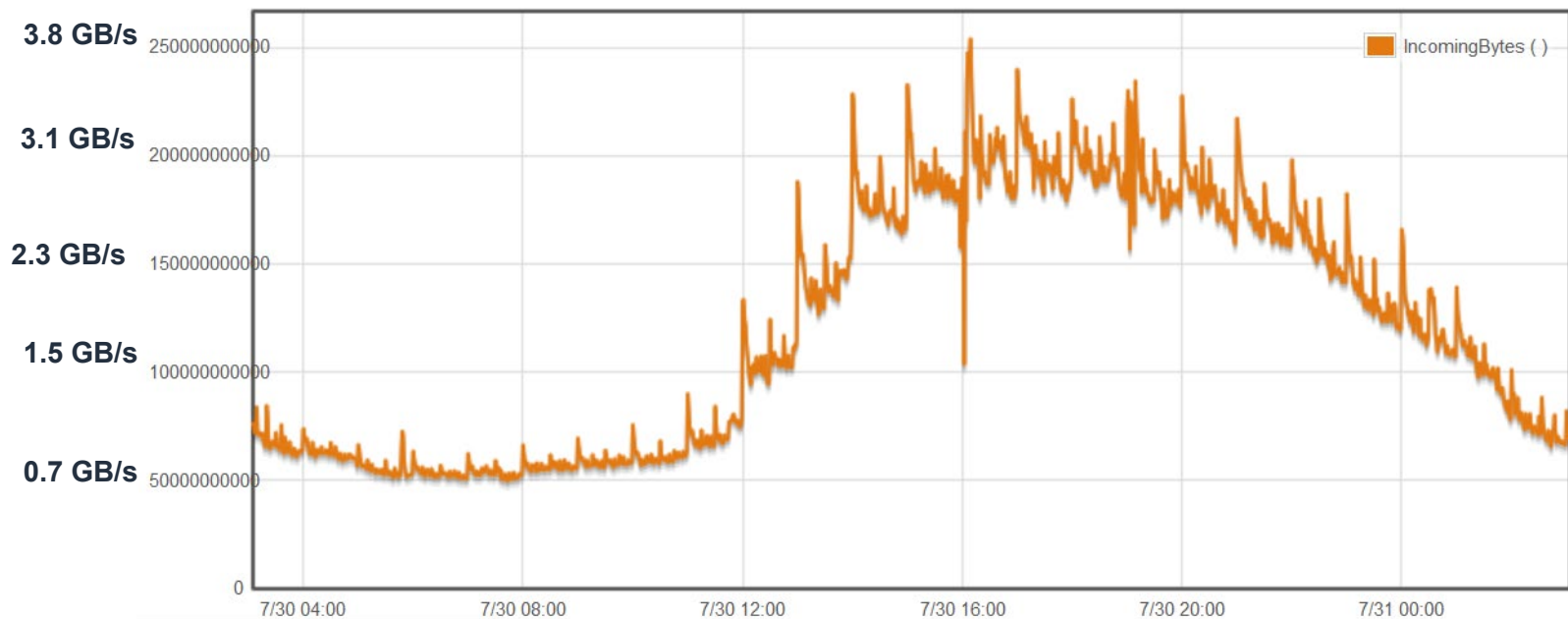


Massive scale & low-latency

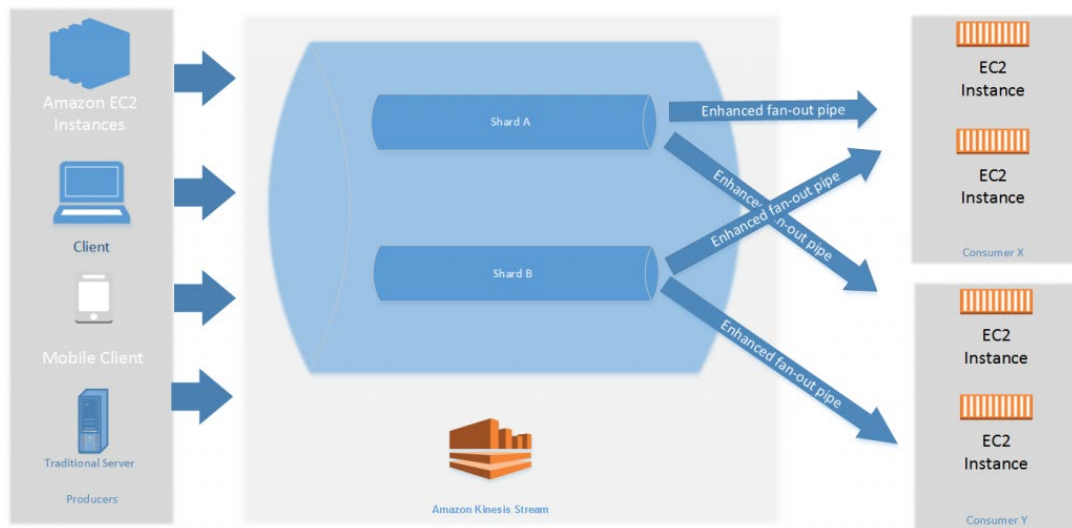


Supports massive scale

You can grow your data stream to support any throughput

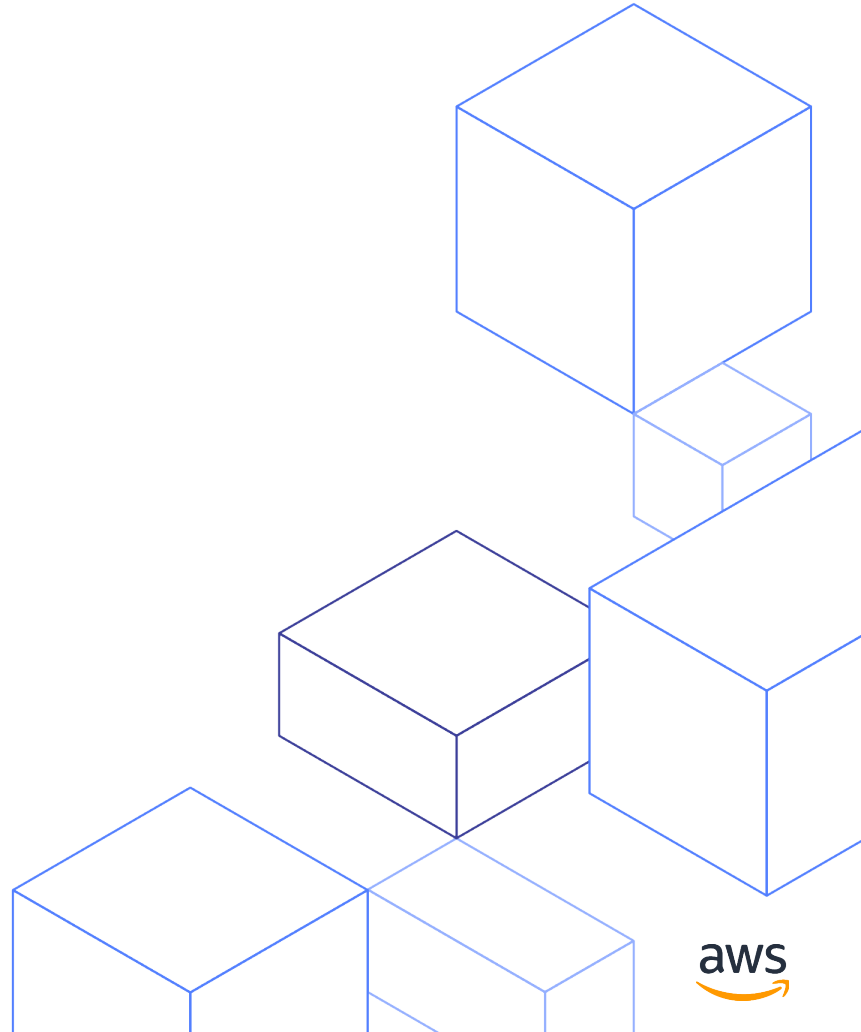


Low-latency and high fan out



- Add 20 consumers
- HTTP/2 to allow <100 ms delivery
- Enhanced Fan Out allows multiple consumers, each at 2MB/second, independently

Low cost



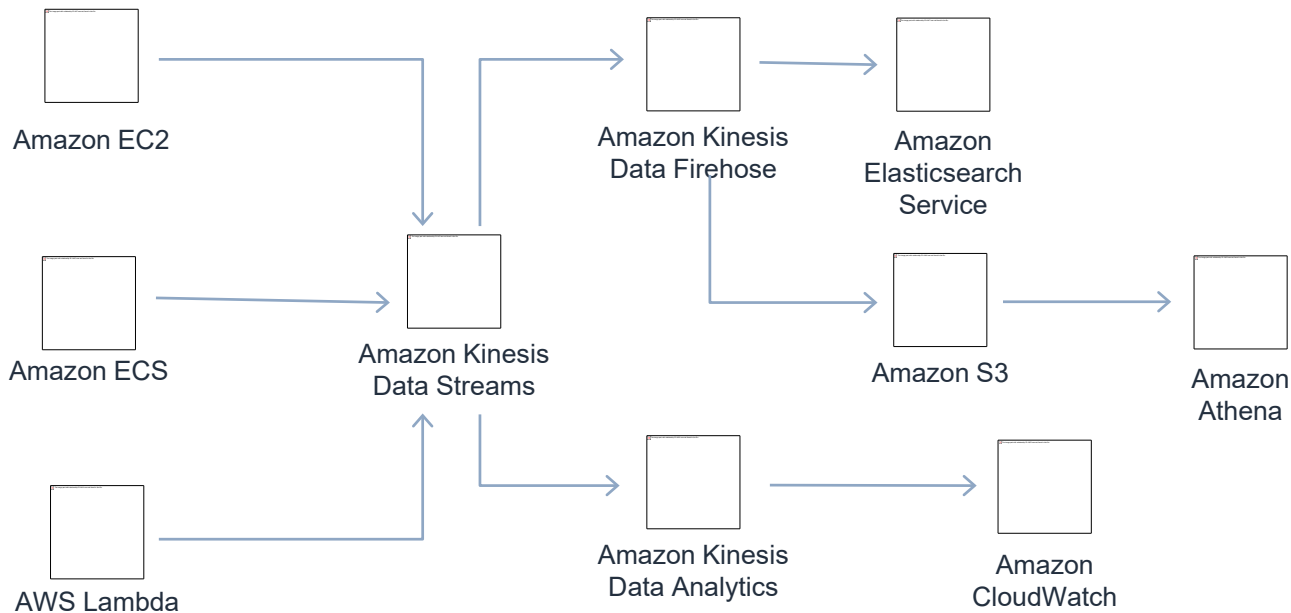
Cost-Effective



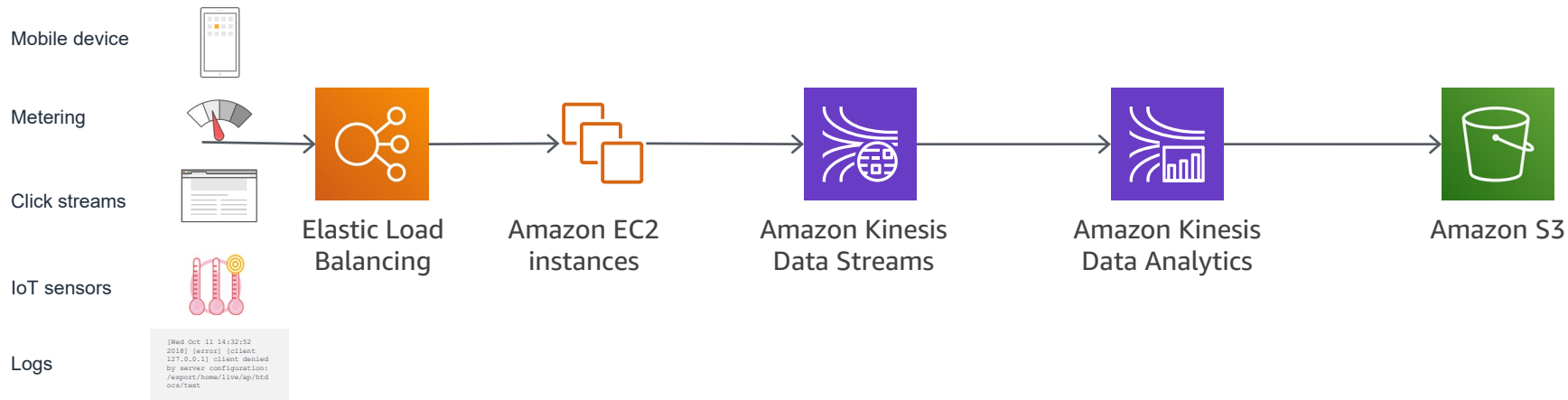
- Pay-as-you-go pricing
- No upfront cost and no minimum fees
- Based on two dimensions:
 - Shard-Hour: \$0.015
 - PUT Payload Units (25K), per million units: \$0.014
- Granular scaling that enables you to balance capacity and costs

Typical data streaming use cases

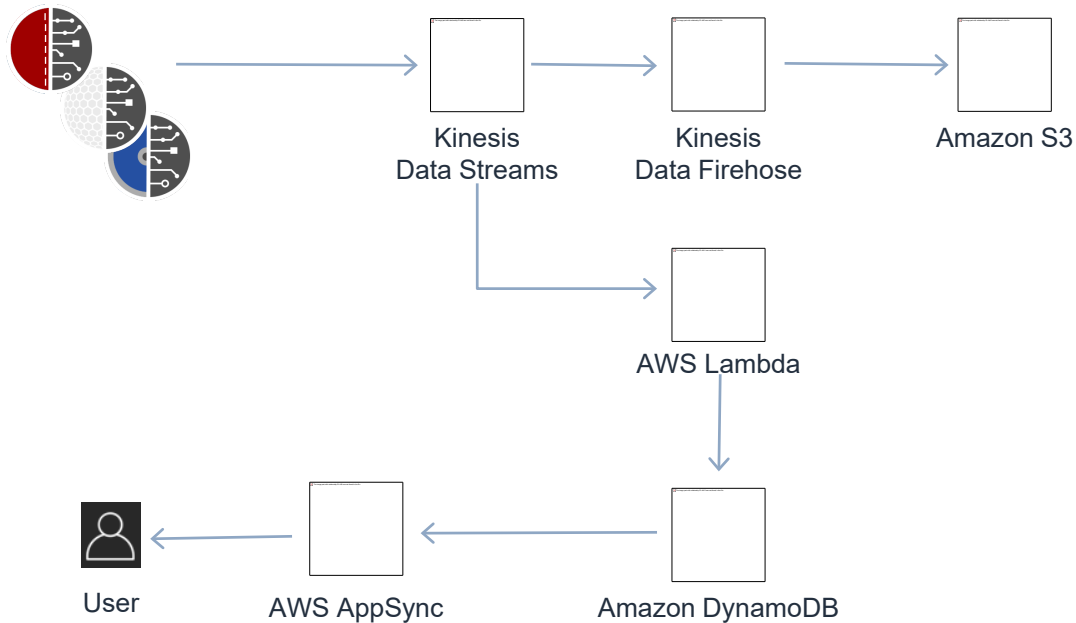
Log ingestion to process Terabytes of data in real time



Streaming ETL to your data lake



IoT sensor data collected, ingested, and analyzed



Recap

- Data streaming opens up possibilities of speed
- Amazon Kinesis makes it easy to build and scale streaming applications at low-cost
- Leverage solution guides and Data Labs to get started

Next steps

Learn more about Amazon Kinesis:

aws.amazon.com/kinesis

Get started with Amazon Kinesis:

aws.amazon.com/kinesis/getting-started