



Overview

Serverless Data Lake Workshop

Objective

- This workshop is meant to give customers a hands-on experience with mentioned AWS services.
- Serverless Data Lake workshop helps customers build a cloud-native and future-proof serverless data lake architecture.
- It allows hands-on time with AWS big data and analytics services including Amazon Kinesis Services for streaming data ingestion and analytics, AWS Glue for ETL and Data Catalogue Management, Amazon Athena to query data lake.

Target Audience

- Data analysts, IT Managers, Solutions Architects, Data Scientists.
- AWS hands-on knowledge: AWS Console, IAM, S3, CloudWatch.
- Concept knowledge: Hadoop, NoSQL, ETL, Kafka Streams, HIVE, Presto engines, JSON, Parquet data formats.

AWS Pop-up Loft | Johannesburg

Audience Check

AWS Pop-up Loft | Johannesburg

Prerequisites

Lab Preparations

- Bring your own AWS Account: You need to use your own AWS Account/user for this workshop. Give AdminRights to the user. (In the lab guide, there are detailed instructions to implement a fine grained policy). It is not recommended to do the workshop from your AWS production account/using root.
- Install the AWS Data Pipeline tool and the KDG tool. Follow the KDG Guide on <https://github.com/awslabs/aws-serverless-data-lake-workshop/tree/main/kdg> to install and configure the KDG on your AWS account. You simply need to use the CloudFormation template to create KDG to us-west-2 (Oregon) region
- Download the Lab Guide & Code: <https://github.com/CloudyYalla/aws-serverless-data-lake-workshop>, including the slides from the workshop

Expected Costs & Cleanup

- You are responsible for the cost of the AWS services used while running this lab.
- As of the date of publication, the baseline cost for running this solution as-is should be around < USD 3.
- Make sure you delete the resources you created.
- The permissions shall either be removed for the user after the immersion day, or they shall be turned into more fine-grained permissions.

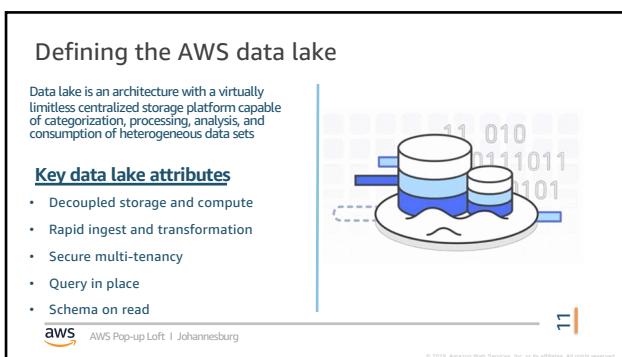
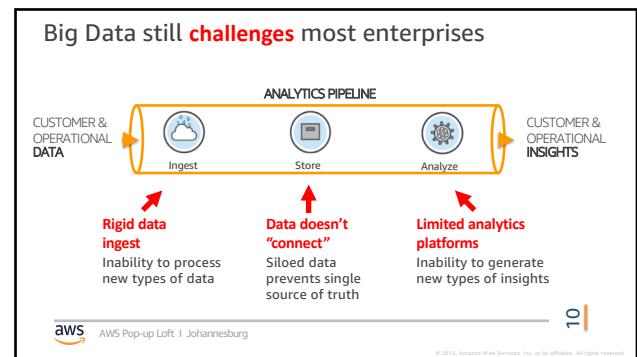
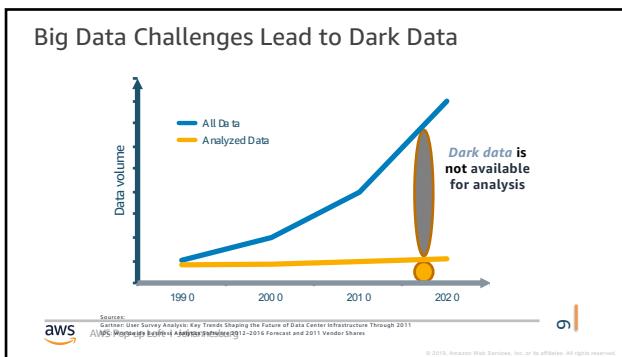
If you don't have AWS account, don't worry, pair with someone or simply watch as go through the hands-on-labs.

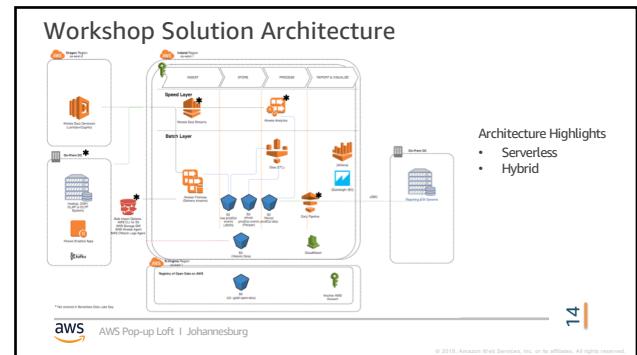
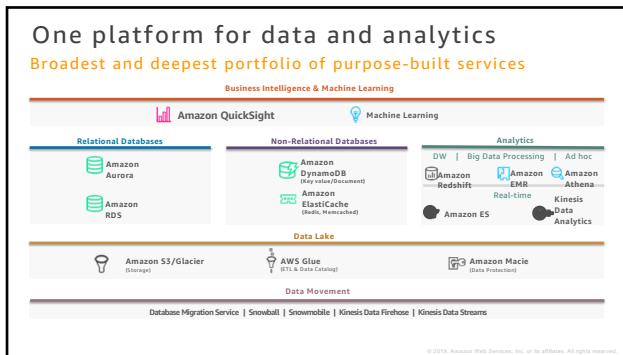
AWS Pop-up Loft | Johannesburg

Hands-on Labs Overview

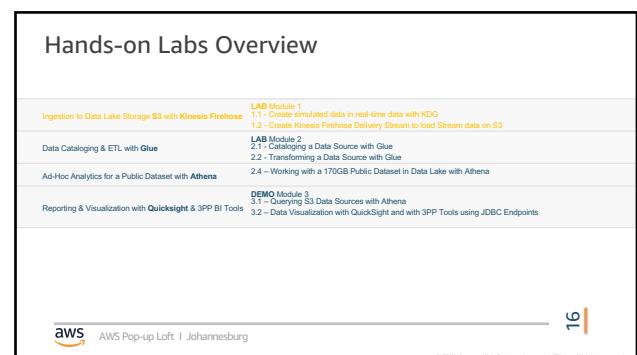
Ingestion to Data Lake Storage S3 with Kinesis Firehose	LAB Module 1
	1.1 – Create Ingested data in real-time data with KDG
	1.2 – Create Kinesis Firehose Delivery Stream to load Stream data on S3
Data Cataloging & ETL with Glue	LAB Module 2
	2.1 – Cataloging a Data Source with Glue
	2.2 – Transforming a Data Source with Glue
Ad-Hoc Analytics for a Public Dataset with Athena	2.4 – Working with a 170GB Public Dataset in Data Lake with Athena
Reporting & Visualization with Quicksight & 3PP BI Tools	DEMO Module 3
	3.1 – Querying S3 Data Sources with Athena
	3.2 – Data Visualization with QuickSight and with 3PP Tools using JDBC Endpoints

AWS Pop-up Loft | Johannesburg





A note on Serverless



Service Used in this Lab

S3: Store data at Any Scalable with Unmatched Durability & Availability

Active Infrequent Archive

AWS Pop-up Loft | Johannesburg

17

- Built to store any amount of data
- Runs on the world's largest global cloud infrastructure
- Designed to deliver 99.99999999% durability
- Geographic redundancy & automatic replication
- Seamlessly replicates data between any region
- Tiered storage to optimize price/performance: Store data at \$0.023/GB/month at S3 (\$0.004/GB/month at Glacier)

Amazon Kinesis—Real Time

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

Kinesis Video Streams	Kinesis Data Streams	Kinesis Data Firehose	Kinesis Data Analytics
Capture, process, and store video streams for analytics	Build custom applications that analyze data streams	Load data streams into AWS data stores	Analyze data streams with SQL

AWS Pop-up Loft | Johannesburg

18

Service Used in this Lab

Amazon Kinesis Firehose

Capture and submit streaming data to Firehose

Firehose loads streaming data continuously into S3, Amazon Redshift, and Amazon ES

Analyze streaming data using your favorite BI tools

Zero administration: Capture and deliver streaming data to Amazon S3, Amazon Redshift, or Amazon Elasticsearch Service without writing an app or managing infrastructure.

Direct-to-data-store integration: Batch, compress, and encrypt streaming data for delivery in as little as 60 seconds.

Seamless elasticity: Seamlessly scales to match data throughput without intervention.

AWS Pop-up Loft | Johannesburg

19

LAB 1

AWS Pop-up Loft | Johannesburg

20

Summary of Lab 1

Congratulations! You have successfully created an ingestion pipeline without servers, which is ready to process huge amounts of data.

In the next lab, you will create the processing pipeline to convert, transform the data from the ingestion layer.

 AWS Pop-up Loft | Johannesburg

21

Hands-on Labs Overview

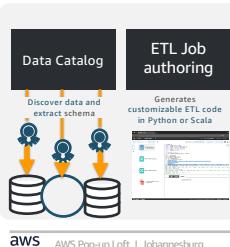
Ingestion to Data Lake Storage S3 with Kinesis Firehose	LAB Module 1 1.1 - Create simulated data in real-time data with KDG 1.2 - Configure Kinesis Firehose Delivery Stream to load Stream data on S3
Data Cataloging & ETL with Glue	LAB Module 2 2.1 - Cataloging a Data Source with Glue 2.2 - Transforming a Data Source with Glue
Ad-Hoc Analytics for a Public Dataset with Athena	2.4 - Working with a 170GB Public Dataset in Data Lake with Athena
Reporting & Visualization with Quicksight & 3PP BI Tools	DEMO Module 3 3.1 - Querying S3 Data Sources with Athena 3.2 - Data Visualization with QuickSight and with 3PP Tools using JDBC Endpoints

 AWS Pop-up Loft | Johannesburg

22

Service used in this Lab

AWS Glue: Data Catalog & ETL Service



23

AWS Glue: components



24

Glue data catalog

Discover and organize your data sets

25

AWS Pop-up Loft | Johannesburg

© 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved.

AWS Glue Data Catalog

Manage table metadata through a Hive metastore API or Hive SQL. Supported by tools such as Hive, Presto, Spark, etc.

We added a few extensions:

- **Search** metadata for data discovery
- **Connection info** – JDBC URLs, credentials
- **Classification** for identifying and parsing files
- **Versioning** of table metadata as schemas evolve and other metadata are updated

Populate using Hive DDL, bulk import, or automatically through **crawlers**.

Table name	Namespace	Location	Classification
Bucket-1	BI Data	S3://MyBucket/MyFold...	Unknown
MarketingExpensesQ42015	BI Data	S3://MyBucket/MyFold...	Unknown
Gl_Schema	BI Data	S3://MyBucket/MyFold...	Clobwatch
MarketingExpensesQ42015_Taxn-0000202015	Default	S3://MyBucket/MyFold...	None
3433344334	Default	S3://MyBucket/MyFold...	Omriture
BrandNewRestHDHC-Table	BI Data	S3://MyBucket/MyFold...	Omriture
3433344334	BI Data	S3://MyBucket/MyFold...	Clobwatch
MyDataset	BI Data	S3://MyBucket/MyFold...	Omriture
Rechrist_Table-00002015	BI Data	S3://MyBucket/MyFold...	Omriture
MyDataset	Default	S3://MyBucket/MyFold...	Web Log
3433344334	Default	S3://MyBucket/MyFold...	Unknown

26

AWS Pop-up Loft | Johannesburg

© 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved.

AWS Glue Data Catalog

A unified metadata repository across data sources & data formats

Integrated with other AWS services

27

AWS Pop-up Loft | Johannesburg

© 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved.

AWS Glue Workflow Overview

Data Stores

1. Use a crawler to infer the schema of your data.
2. Use Classifiers to Populate the AWS Glue Data Catalog with table definitions.
3. Crawler connects to the data store.
4. Schema inference occurs.
5. Define a job that describes the transformation of data from source to target.
6. Run your job to transform your data.
7. Monitor

28

AWS Pop-up Loft | Johannesburg

© 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved.

AWS Glue Built-in Classifiers

Ref: AWS Glue Developer Guide

29

Crawlers: Auto-populate Data Datalog

Automatic schema inference:

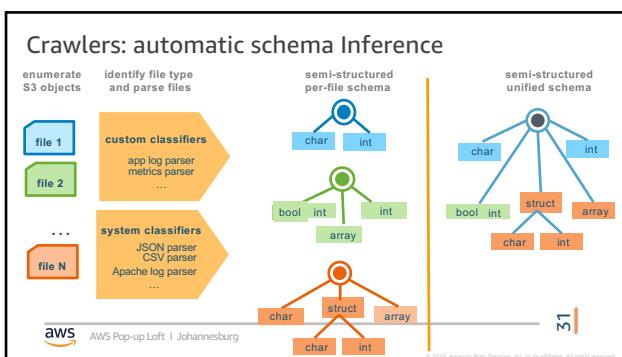
- Built-in classifiers detect file type and extract schema: record structure and data types.
- Add your own or share with others in the Glue community - It's all Grok and Python.

Auto-detects Hive-style partitions, grouping similar files into one table.

Run crawlers on schedule to discover new data and schema changes.

Serverless – only pay when crawls run.

30

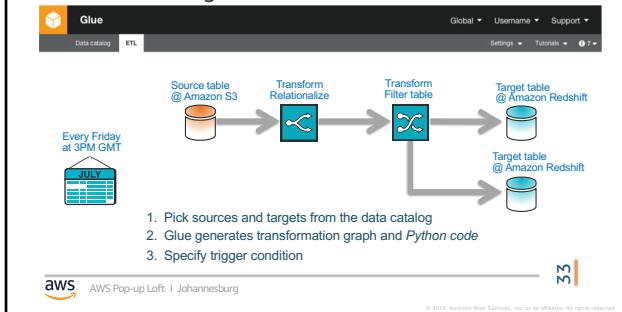


Job authoring in Glue

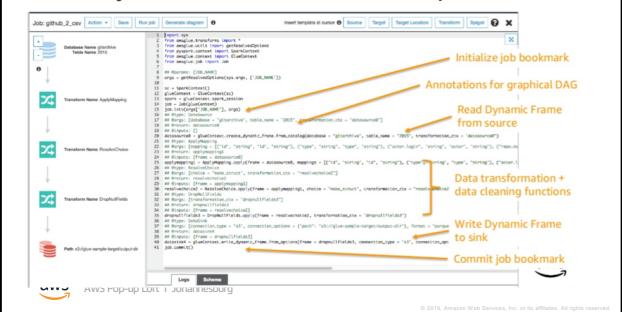
Make ETL job authoring like code development using your own tools

32

Automatic code generation



Anatomy of a Generated ETL Job Script

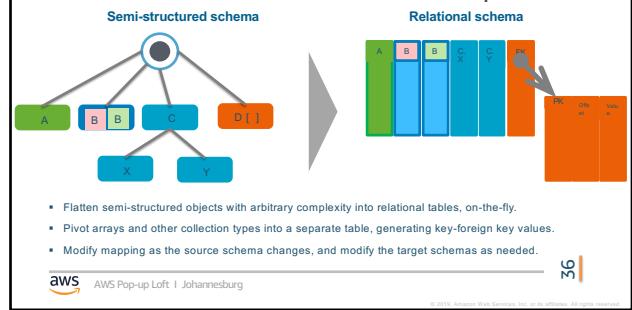


Scheduling Jobs

The screenshot shows the AWS Lambda 'Add trigger' interface with two tabs:

- Several event types**: Shows a list of triggers including CloudWatch Metrics, CloudWatch Logs, CloudWatch Events, S3, CloudWatch Metrics Insights, and CloudWatch Metrics Insights Metrics.
- Pass parameters**: Shows a list of triggers including CloudWatch Metrics, CloudWatch Logs, CloudWatch Events, S3, CloudWatch Metrics Insights, CloudWatch Metrics Insights Metrics, and Lambda functions.

Glue transformations are flexible and adaptive



Kahoot!

AWS Pop-up Loft | Johannesburg

37

Edit & Test with Dev-Endpoints

Connect your IDE to an AWS Glue development endpoint Environment to interactively develop, debug & test ETL code.

AWS Pop-up Loft | Johannesburg

38

Explorer and Experiment with Data

Connect your notebook (e.g. Jupyter, Zeppelin) to an AWS Glue development endpoint.

Interactively experiment and explore datasets and data sources.

Deploy to production
Push scripts to S3
Create or register with ETL job

AWS Pop-up Loft

39

Orchestration & resource management

Fully managed, serverless job execution

AWS Pop-up Loft | Johannesburg

40

Serverless Job Execution

There is no need to provision, configure, or manage servers

- Warm pools: pre-configured fleets of instances to reduce job startup time
- Auto-configure VPC and role-based access
- Automatically scale resources to meet SLA and cost objectives
- Only pay for the resources you consume per-second billing (10-minute min)

41

© 2018 Amazon Web Services, Inc. or its Affiliates. All rights reserved.

AWS Pop-up Loft | Johannesburg

Apache Spark & AWS Glue ETL

What is Apache Spark?

- Parallel, scale-out data processing engine
- Fault-tolerance built-in
- Flexible interface: Python scripting, SQL
- Rich eco-system: ML, Graph, analytics, ...

Apache Glue ETL libraries

- Integration: Data Catalog, job orchestration, code-generation, job bookmarks, S3, RDS
- ETL transforms, more connectors & formats
- New data structure: Dynamic frames

42

© 2018 Amazon Web Services, Inc. or its Affiliates. All rights reserved.

AWS Pop-up Loft | Johannesburg

Performance: Lots of small files

AWS Glue ETL small file scalability

1.2 Million Files

Time (sec)

partitions : # files

43

© 2018 Amazon Web Services, Inc. or its Affiliates. All rights reserved.

AWS Pop-up Loft | Johannesburg

Dynamic Frame Transformations

15+ transforms out-of-the box

44

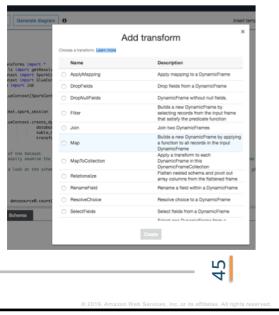
© 2018 Amazon Web Services, Inc. or its Affiliates. All rights reserved.

AWS Pop-up Loft | Johannesburg

Useful AWS Glue Transforms

- **toDF():** Convert to a Dataframe
- **Spigot():** Sample data of any Dynamic Frame to S3
- **Filter(), map():** Apply python UDFs to Dynamic Frames
- **Join():** Join two Dynamic Frames

and more...



45

aws AWS Pop-up Loft | Johannesburg

Hands-on Labs Overview

Ingestion to Data Lake Storage S3 with Kinesis Firehose	LAB Module 1 1.1 - Ingest unstructured data in real-time with KDG 1.2 - Create Kinesis Firehose Delivery Stream to load Stream data on S3
Data Cataloging & ETL with Glue	LAB Module 2 2.1 - Cataloging a Data Source with Glue 2.2 - Transforming a Data Source with Glue
Ad-Hoc Analytics for a Public Dataset with Athena	2.4 - Working with a 170GB Public Dataset in Data Lake with Athena
Reporting & Visualization with Quicksight & 3PP BI Tools	DEMO Module 3 3.2 - Data Visualization with QuickSight and with 3PP Tools using JDBC Endpoints

46

aws AWS Pop-up Loft | Johannesburg

© 2018 Amazon Web Services, Inc. or its Affiliates. All Rights Reserved.

LAB 2.4

aws AWS Pop-up Loft | Johannesburg

47

Service used in this Lab

Amazon Athena—Interactive Analysis

Interactive query service to analyze data in Amazon S3 using standard SQL
No infrastructure to set up or manage and no data to load
Ability to run SQL queries on data archived in Amazon Glacier (coming soon)

Query Instantly	Pay per query	Open	Easy

Zero setup cost: just point to S3 and start querying

Pay only for queries run; save 30–90% on per-query costs through compression

ANSI SQL Interface, JDBC/ODBC drivers, input formats, compression types, and complex joins and data types

Serverless: zero infrastructure, zero administration. Integrated with QuickSight

48

aws AWS Pop-up Loft | Johannesburg

Public Dataset

The [Global Database of Events, Language and Tone \(GDELT\) Project](http://www.gdeltproject.org) monitors the world's broadcast, print, and web news from nearly every corner of every country in over 100 languages and identifies the people, locations, organisations, counts, themes, sources, emotions, counts, quotes, images and events driving our global society every second of every day. The data set v1.0 is publicly available in S3 in the [Registry of Open Data on AWS](#).

AWS Pop-up Loft | Johannesburg

Ad-Hoc Query from the Data Lake:

Finding the aggregate number of events for the last 25+ years

Notice the data amount scanned? The results are returned by scanning 170+ GB of data from 4000+ uncompressed CSV files on S3. That's the power of HIVE, Presto and other Hadoop Technologies simplified by Athena Service.

AWS Pop-up Loft | Johannesburg

HIVE DDL Commands

- We are creating a schema definition in our Glue service, in our Data Catalogue.
- The actual data is in another AWS account, and in another AWS region (us-east-1: Northern Virginia)
- You can Access this data, because it is a public dataset located in 's3://gdelt-open-data/events/folder', and is open to everyone.
- Although we are creating TABLEs, there is no database. Events table is a representation of thousands of TSV (Tab Separated Files) files stored in S3.
- Technologies like Apache HIVE and Presto enables accessing them using SQL like expressions.

51

AWS Pop-up Loft | Johannesburg

Summary of Lab 2

Congratulations. You have successfully catalogued your data, created serverless ETL jobs for transforming your data.

You have also performed an ad-hoc query of thousands of uncompressed CSV files (containing hundreds of millions of lines, more than 170 GB) from data lake storage using Athena.

52

AWS Pop-up Loft | Johannesburg

Hands-on Labs Overview

Ingestion to Data Lake Storage S3 with Kinesis Firehose LAB Module 1
1.1 - Create unstructured data in real-time data with KOG
1.2 - Create a Kinesis Firehose Delivery Stream to load Stream data on S3

Data Cataloging & ETL with Glue LAB Module 2
2.1 - Cataloging a Data Source with Glue
2.2 - Transforming a Data Source with Glue

Ad-Hoc Analytics for a Public Dataset with Athena LAB Module 3
2.4 - Working with a 170GB Public Dataset in Data Lake with Athena

Reporting & Visualization with Quicksight & 3PP BI Tools DEMO Module 3
3.2 - Data Visualization with QuickSight and with 3PP Tools using JDBC Endpoints

AWS Pop-up Loft | Johannesburg

53

Demo

*Due to time constraint, we will make a demo of Lab 3.
You are encouraged to do it yourself after the session using the Lab Guide documents.

AWS Pop-up Loft | Johannesburg

54

Amazon QuickSight

Amazon QuickSight is a fast, cloud-powered business intelligence (BI) service that makes it easy for you to deliver insights to everyone in your organization.

Data Scientist (Author)
Give power users and analysts the freedom to do their own self-serve data discovery and analysis on governed data you control

Dashboard Creator (Author)
Create and publish rich, interactive dashboards to all of your users

End User (Reader)
With the new Reader Role, you can provide everyone in your organization secure, easy access to interactive dashboards and reports, on any device

Why Amazon QuickSight?

No servers to manage
Amazon QuickSight has no servers or software to manage, maintain, deploy, upgrade or migrate. We do the heavy lifting so you don't have to.

Easily scale from 10 users to 10,000
QuickSight automatically scales with your usage and cost with no need for manual administration. QuickSight will grow with your organization's needs: from a few users to tens of thousands of users.

Native AWS integration
Amazon QuickSight securely integrates with your data sources and AWS services like Amazon Simple Storage Service (Amazon S3), Redshift, Amazon Athena, Amazon Aurora, Amazon Relational Database Service (Amazon RDS), Amazon CloudWatch Metrics, Amazon Identity and Access Management (IAM), AWS CloudTrail, Amazon Cloud Directory and more - providing you with everything you need to build an end-to-end BI solution.

Pay only for what you use
Power read-only access to interactive dashboards and pay only when your users access them with Pay-per-Session pricing. With Amazon QuickSight there are no upfront costs, no annual commitments and no charges for inactive users.

You are Invited to the next Loft Session about QuickSight on Monday

Serverless BI Using Amazon QuickSight

11th of March 2019
14:00-15:00 @Johannesburg



AWS Pop-up Loft | Johannesburg

Kahoot! time

AWS Pop-up Loft | Johannesburg

58

Summary

AWS helps unlock Big Data in three ways



- Future-proofed for new data types**
Ingest new data for differentiated experiences: image/video, social media, IoT sensors, and more
- Single source of truth (aka "data lake")**
Secure and scalable backbone for all incoming data and all analytics – no more data siloes
- The right analytics tool for every job**
New analytics possibilities – machine learning, natural language processing, Hadoop-as-a-service, and more

59

AWS Pop-up Loft | Johannesburg

Summary

Benefits of Data Lakes from AWS



- Open and comprehensive
- Secure
- Scalable and durable
- Lowest cost

AWS Pop-up Loft | Johannesburg

60

