FSV303

# AWS re:INVENT

## Queryable Archive + Data Lake

FSV303 - Building Queryable Archives and Data Lakes for Financial Services

George Smith, Global Financial Services Solutions Architect at AWS

November 27, 2017

AWS re:Invent

aws

Financial institutions today must manage multiple data types from a wide variety of sources. Among these various data types, archive data presents a particular challenge: it is invisible to much of the organization and not easily leveraged by the lines of business for analytics, insight, and product innovation. Faced with massive volumes of archive data, Financial Services organizations are finding that delivering insights in a timely manner requires a data storage and analytics solution with more agility and flexibility than traditional data management systems can provide. In this session, we will discuss a design pattern that (1) brings this data into a highly available, lower-cost queryable archive within AWS than you currently have and (2) migrates that data to a data lake that the entire organization can use to extract insight and drive innovation. We will walk through a strategy that addresses the following topics: storing archive data in compressed, cost-effective, and readily available formats; creating lifecycle policies to archive older data sets and make them easily accessible; fully utilizing the features of object storage to enrich the data lake; and applying AWS analytics tools to gather business insights.
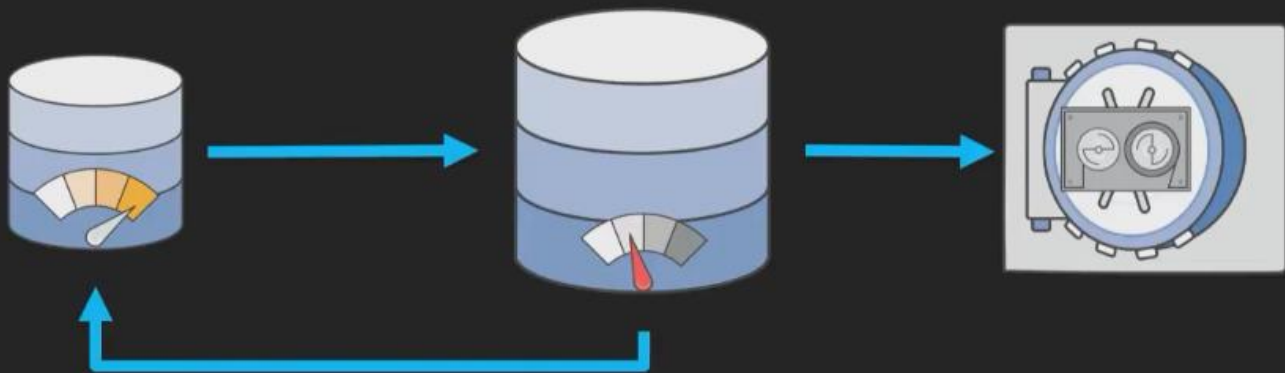
## What to Expect From This Session

- A pattern for better, cheaper, faster archives called "Queryable Archive"
- An implementation of "Queryable Archive + Data Lake" on AWS
- Why bringing archived data online exposes "dark data"

# Agenda

- Why We Archive
- Building A "Queryable Archive + Data Lake"
- Demo of A "Queryable Archive + Data Lake"
- Benefits and Costs
- Unleashing "Dark Data"
- Next Steps

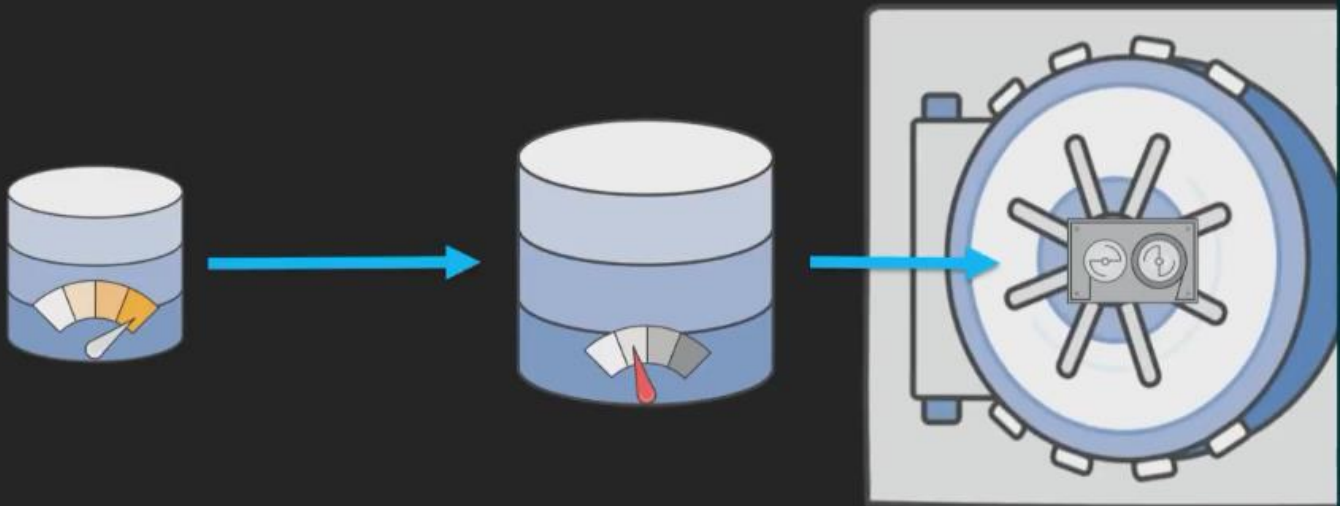# Why We Archive

# Archive For Performance

Archive For Regulatory

We now have way more storage attached to compute which is not what we intended during design phase
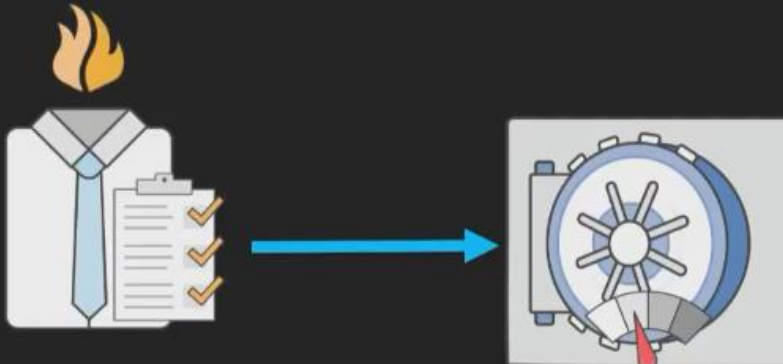


Access For Regulatory, Legal, Customer, etc.

we now have an end user who needs access to data and can't get it



We have to set up spare database and copy the archives into it, then get the date out for the customer

We will then tear everything down again

Building A "Queryable Archive + Data Lake"*

* (Without breaking the bank!)


A Queryable Archive + Data Lake

Extract    Transform    Load    Explore

We start out with all our historical financial datasets that are sitting in our databases in data centers today. Once the data is loaded into the specific AWS database services, we have a variety of tools to start exploring that data


Queryable Archive + Data Lake Tech Stack

| | |
|---|---|
| Amazon S3 (Simple Storage Service) | Secure, durable, highly scalable cloud storage |
| AWS Database Migration Service | Helps you migrate your database to AWS |
| AWS Glue | Fully Managed ETL Service |
| Amazon Redshift | Petabyte scale Data Warehouse solution |
| Amazon Redshift Spectrum | Redshift SQL queries against exabytes in Amazon S3 |
| Amazon EMR | Managed Hadoop framework |
| Amazon Athena | Serverless interactive query service |

We use the AWS Data Migration Service and now have the data in S3 in a compressed CSV format.



We now again use AWS Glue to transform the data into our enterprise data format and store the data in S3 in the Parquet object format

We can now access and explore the data

A Queryable Archive + Data Lake On AWS

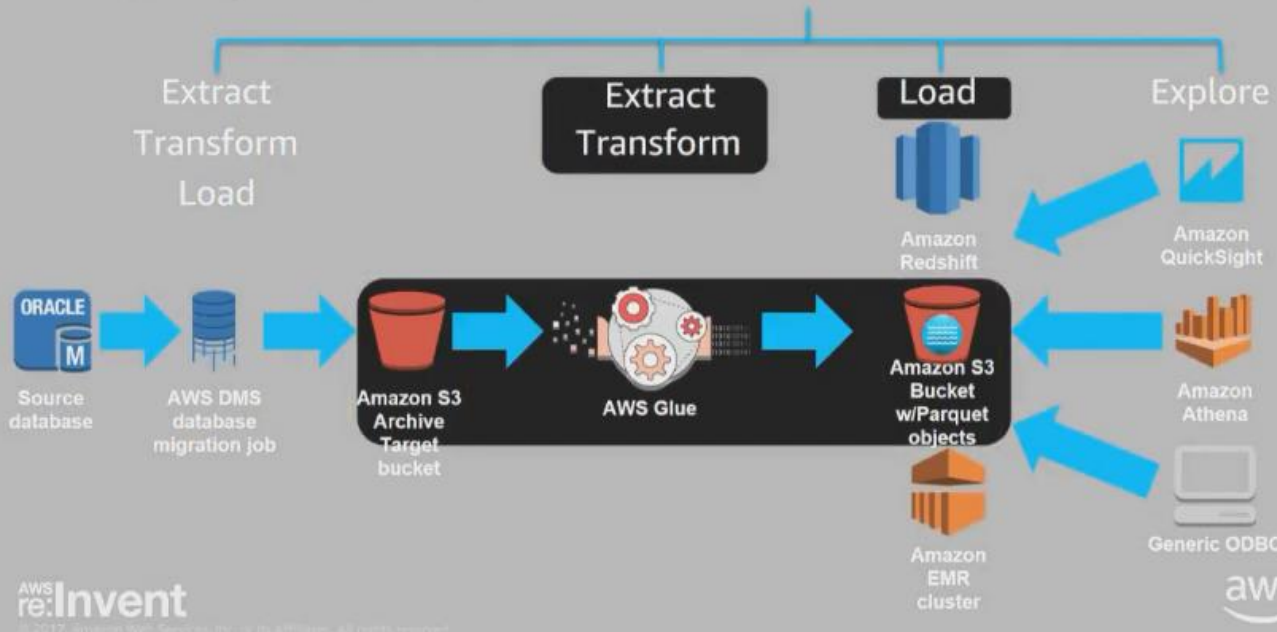A Task is a batch job that moves data from one point to another

The job is done 4h 26m later and we have copied 167,497,743 rows from the Oracle database to S3 bucket in compressed CSV format. The dataset contains orders and order executions by a trading desk.
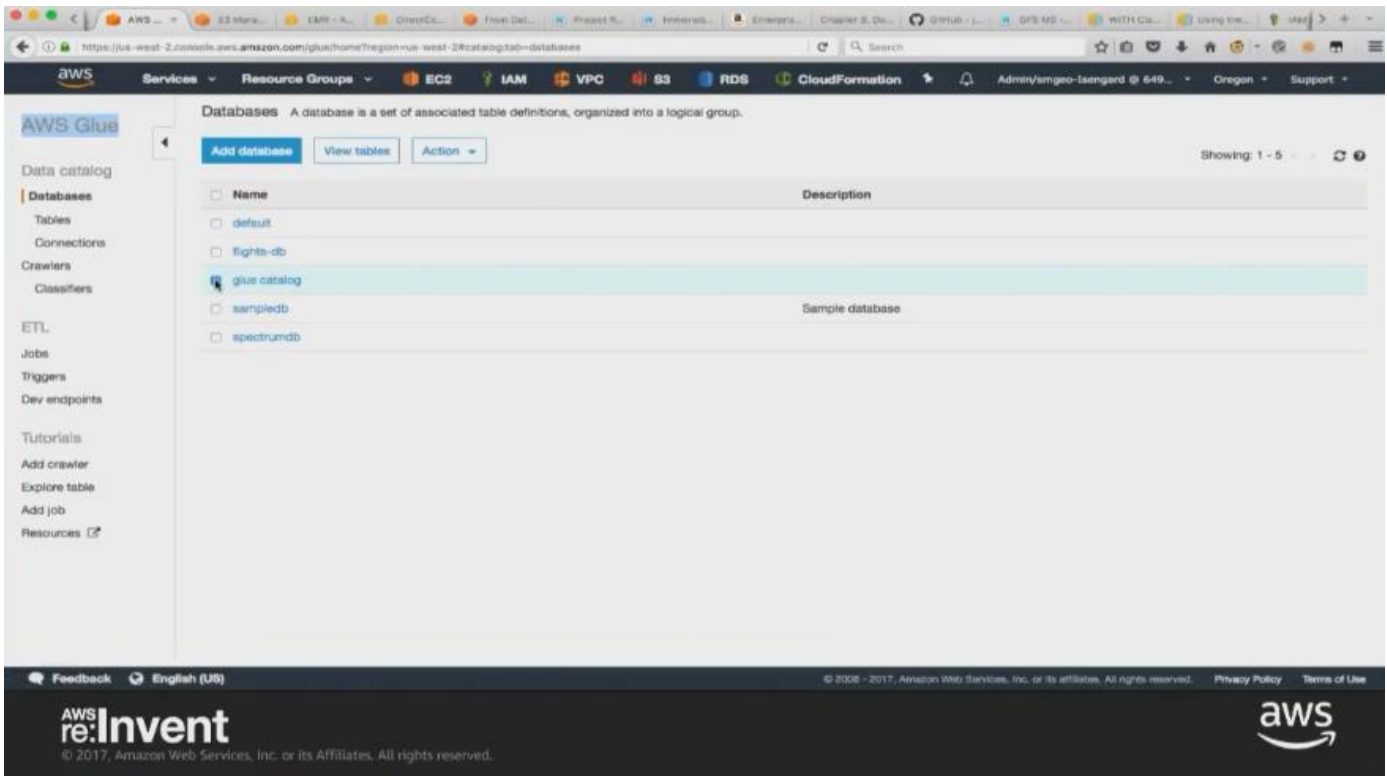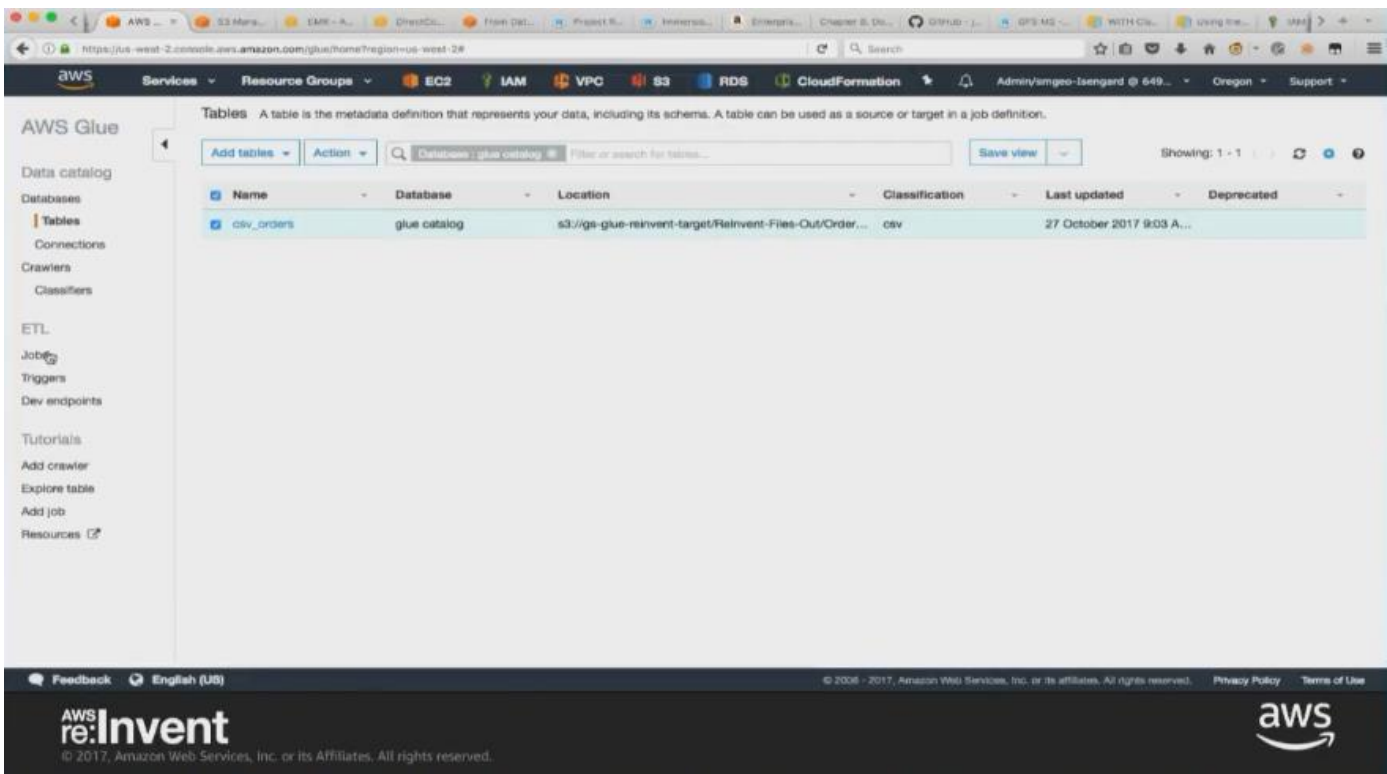
A Queryable Archive + Data Lake On AWS



A Queryable Archive + Data Lake On AWS

We now want to transform the compressed CSV data into the enterprise data lake format for the data model in Parquet objects. You will be transforming and enriching the data along the way

We are going to be using AWS Glue for this batch job,



Within the glue_catalog database, we have the table containing our compressed CSV files

AWS Glue generate python code for the CSV_to_Parquet job for us

The job finishes after 2h 55mins and we have loaded all 167,497,743 rows into our data lake.

Now we have data in 2 locations.

We also load data into Redshift

```
f45c89bd7fe5:~ smgeo$ aws emr create-cluster \
> --termination-protected \
> --applications Name=Hadoop Name=Hive Name=Pig Name=Hue Name=Spark Name=Presto \
> --tags 'name=queryarchive' \
> --ec2-attributes '{"KeyName":"MyEC2Key","InstanceProfile":"EMR_EC2_DefaultRole","SubnetId":"subnet-17b0fa73","EmrManagedSlaveSecurityGroup":"sg-ea673093","EmrManagedMasterSe
curityGroup":"sg-eb673092"}' \
> --service-role EMR_DefaultRole \
> --enable-debugging \
> --release-label emr-5.9.0 \
> --log-uri 's3n://aws-logs-649225637812-us-west-2/elasticmapreduce/' \
> --name 'QueryArchive Cluster Glue' \
> --instance-groups '[{"InstanceCount":1,"InstanceGroupType":"MASTER","InstanceType":"m3.xlarge","Name":"Master - 1"},{"InstanceCount":2,"InstanceGroupType":"CORE","InstanceTy
pe":"r3.2xlarge","Name":"Core - 2"}]' \
> --region us-west-2 \
> --configurations file://glueConfiguration.json
{
    "ClusterId": "j-3JJX71R6GHGVP"
}
f45c89bd7fe5:~ smgeo$
```



```
f45c89bd7fe5:~ smgeo$ aws emr create-cluster \
> --termination-protected \
> --applications Name=Hadoop Name=Hive Name=Pig Name=Hue Name=Spark Name=Presto \
> --tags 'name=queryarchive' \
> --ec2-attributes '{"KeyName":"MyEC2Key","InstanceProfile":"EMR_EC2_DefaultRole","SubnetId":"subnet-17b0fa73","EmrManagedSlaveSecurityGroup":"sg-ea673093","EmrManagedMasterSe
curityGroup":"sg-eb673092"}' \
> --service-role EMR_DefaultRole \
> --enable-debugging \
> --release-label emr-5.9.0 \
> --log-uri 's3n://aws-logs-649225637812-us-west-2/elasticmapreduce/' \
> --name 'QueryArchive Cluster Glue' \
> --instance-groups '[{"InstanceCount":1,"InstanceGroupType":"MASTER","InstanceType":"m3.xlarge","Name":"Master - 1"},{"InstanceCount":2,"InstanceGroupType":"CORE","InstanceTy
pe":"r3.2xlarge","Name":"Core - 2"}]' \
> --region us-west-2 \
> --configurations file://glueConfiguration.json
{
    "ClusterId": "j-3JJX71R6GHGVP"
}
f45c89bd7fe5:~ smgeo$ more glueConfiguration.json
[
    {
    "Classification": "hive-site",
    "Properties": {
        "hive.metastore.client.factory.class": "com.amazonaws.glue.catalog.metastore.AWSGlueDataCatalogHiveClientFactory"
    }
    },
    {
    "Classification": "spark-hive-site",
    "Properties": {
        "hive.metastore.client.factory.class": "com.amazonaws.glue.catalog.metastore.AWSGlueDataCatalogHiveClientFactory"
    }
    },
    {
    "Classification": "spark-log4j",
    "Properties": {
        "log4j.rootCategory": "WARN, console"
        }
    }
]

f45c89bd7fe5:~ smgeo$
```

# Benefits

- Single golden source of truth
- Right sizing compute for workload
- **Increased availability**

# A Queryable Archive + Data Lake On AWS

# Storage Efficiencies



Bar chart titled "Storage Efficiencies" with y-axis labeled GB (0.00 to 80.00):
- Oracle: ~75
- Amazon S3 CSV Uncompressed: ~56
- Amazon Redshift: ~17
- Amazon S3 CSV Compressed: ~5
- Amazon S3 Parquet: ~3

# Compute Efficiencies

| Storage Format | Size in GB | Size as a percent of largest store* | Costs Per Month |
|---|---|---|---|
| RDS Oracle | 75.70 | 100.00% | $1449.78 / $706.69 |
| Amazon S3 CSV Uncompressed | 56.40 | 74.50% | $1.30 |
| Amazon Redshift | 17.30 | 22.85% | $366.00 |
| Amazon S3 CSV Compressed | 5.43 | 7.17% | $0.13 |
| Amazon S3 Parquet | 3.40 | 4.49% | $0.08 |

## Data Efficiencies

| Storage Tier | Monthly Storage Pricing per PB |
|---|---|
| Amazon S3 | $22,583.30 |
| Amazon S3-IA | $13,107.20 |
| Amazon Glacier | $4,194.31 |

## Unleashing "Dark Data"

## What is Dark Data?

"Buried within raw information generated in mind-boggling volumes by transactional systems . . . are critical strategic, customer, and operational insights that, once illuminated by analytics, can validate or clarify assumptions, inform decision making, and help chart new paths to the future."

Kambies, T., Roma, P., Mittal, N., & Sharma, S. K. (2017, February 7). Dark analytics: Illuminating opportunities hidden within unstructured data. Retrieved October 16, 2017, from https://dupress.deloitte.com/dup-us-en/focus/tech-trends/2017/dark-data-analyzing-unstructured-data.html

A Queryable Archive + Data Lake On AWS



A Flywheel For Data

Next Steps

## What to Expect From This Session

✓ • A pattern for better, cheaper, faster archives called "Queryable Archive"

✓ • An implementation of "Queryable Archive + Data Lake" on AWS

✓ • Why bringing archived data online exposes "dark data"

## What you can do next?

• Go back and determine the cost to your company of your archive data stores

• What would it cost if you moved them into AWS with this solution?
   • Estimate AWS cost with the simple monthly calculator:
     https://calculator.s3.amazonaws.com/index.html

• Ask yourself, can this increase my organization's efficiency?

• Go build it!

aws

AWS re:Invent

## Thank you

George Smith, Global Financial Services Solutions Architect at AWS
smgeo@amazon.com

aws