



## Data Pipelines with AWS Glue (Level 200)

Unni Pillai, Specialist Solution Architect Thanomsak Ajjanapanya, Data Engineering Manager, Toyota Tsusho – Thailand

## In this session...

Introduction to AWS Glue

**Understand AWS Glue Components** 

Construct an ETL Flow

Demo – Build a data pipeline using Glue in 4 steps

Learn how NEXTY Electronics is using AWS Glue





## What is AWS Glue?

# Fully-managed, serverless extract-transform-load (ETL) service for developers, built by developers

1000s of Developers and jobs





## There are many tools already in AWS Ecosystem

#### **Amazon Redshift Partner Page for Data Integration**





































## Still ETL Developers Hand-Code

Canvas based tools are hard to extend

- Code is flexible, powerful, and easy to share
- Familiar tools and development pipelines
  - IDEs, version control, testing, continuous integration
- Highly customizable tasks can be achieved





## Hand-coding is laborious

schemas change
data formats change
add or change sources
data volume grows

makes hand-coding error-prone & brittle

AWS Glue does the undifferentiated heavy lifting so developers can easily customize





## **AWS Glue Components**



**Data Catalog** 

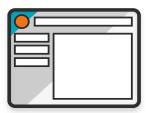
**Discover** 

Automatic crawling

Apache Hive Metastore compatible

Integrated with AWS analytic

services



**Job Authoring** 

Develop

Auto-generates ETL code
Python and Apache Spark
Edit, Debug, and Explore



Serverless execution

Flexible scheduling

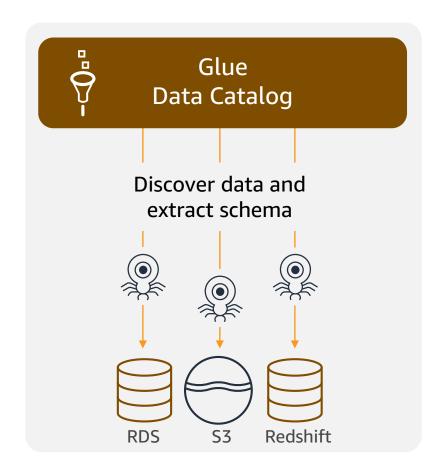
Monitoring and alerting





## AWS Glue - data catalog

Make data discoverable



Automatically discovers data and stores schema

Catalog makes data searchable, and available for ETL

Catalog contains table and job definitions

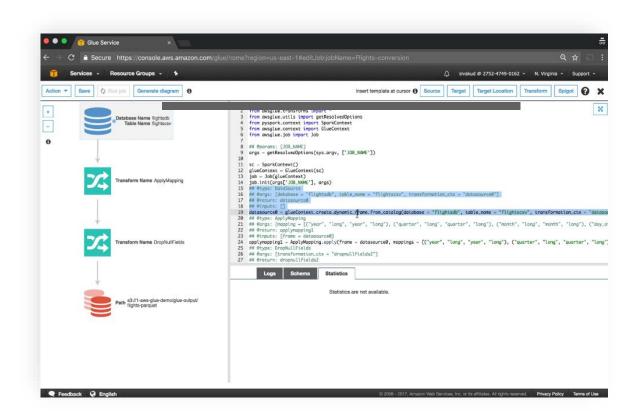
Computes statistics to make queries efficient





## AWS Glue - ETL service

Make ETL scripting and deployment easy



#### **Serverless** Transformations

Based on Apache Spark

Automatically generates ETL code

Code is customizable with PySpark and Scala

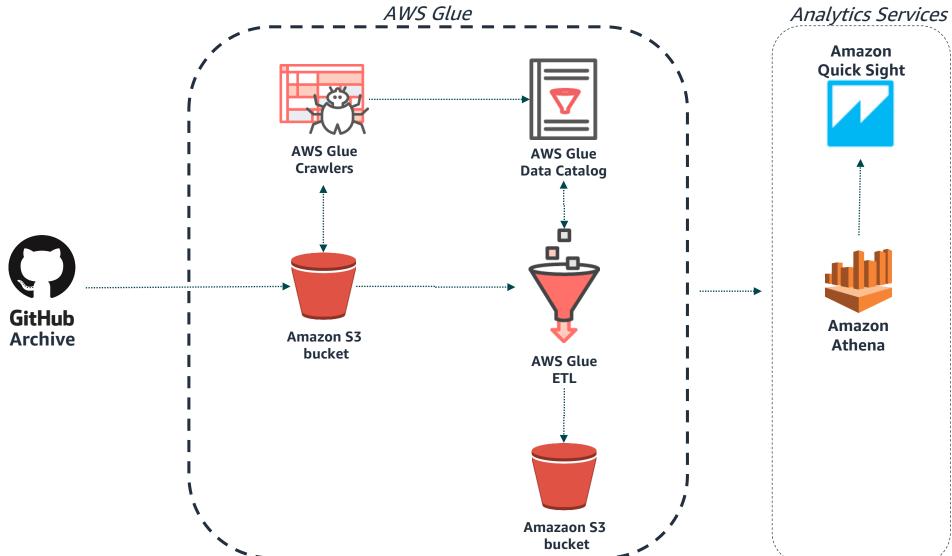
Endpoints provided to edit, debug, test code

Jobs are scheduled or event-based





## ETL example







## Apache Spark and AWS Glue ETL



SparkSQLAWS Glue ETLDataframesDynamic FramesSpark core: RDDs

#### 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, ...

#### **AWS Glue ETL libraries**

Integration: Data Catalog, job orchestration, code-generation, job bookmarks, S3, RDS

ETL transforms, more connectors & formats

New data structure: Dynamic Frames





## Public GitHub timeline is ....

```
Terminal - emacs-25.2 - 130×42
 ile Edit Options Buffers Tools Javascript Help
 'id":"2489651045",<mark>"type":"CreateEvent"</mark>,"actor":{"id":<mark>665991</mark>,"login":"petroav","gravatar_id":"","url":"https://api.github.com/use
rs/petroav", "avatar_url": "https://avatars.githubusercontent.com/u/665991?"}, "repo": {"id": 28688495, "name": "petroav/6.828", "url": "h
ttps://api.github.com/repos/petroav/6.828"},"payload":{"ref":"master","ref type":"branch","master branch":"master","description"
      'public":true."created at":"2015-01-01T15:00:00Z"
"id":"2489651051",<mark>"type":"PushEvent"</mark>,"actor":{"id":3854017,"login":"rspt","gravatar_id":"","url":"https://api.github.com/users/r
    "avatar url":"https://avatars.githubusercontent.com/u/3854017?"},"repo":{"id":28671719,"name":"rspt/rspt-theme","url":"https,
//api.github.com/repos/rspt/rspt-theme"}, "payload": {"push id":536863970, "size":1, "distinct size":1,
              <mark>acfcd"}]},</mark>"public":true,"created at":"2015-01-01T15:00:01Z"}
 "id":"2489651057"<mark>,"type":"WatchEvent"</mark>,"actor":{"id":6894991,"login":"SametSisartenep","gravatar_id":"","url":"https://api.github
.com/users/Samet$isartenep","avatar_url":"https://avatars.githubusercontent.com/u/6894991?"},"repo":{"id":2871998,"name":"visionm
edia/debug", "url": "https://api.github.com/repos/visionmedia/debug"}, "payload": {"action": "started"}, "public": true, "created at": "20
15-01-01T15:00:03Z","org":{"id":9285252,"login":"visionmedia","gravatar_id":"","url":"https://api.github.com/orgs/visionmedia","a
vatar url":"https://avatars.githubusercontent.com/u/9285252?"}}
 "id":"2489651091",<mark>"type":"IssuesEvent"</mark>,"actor":{"id":6269456,"login":"yhoonkim","gravatar_id":"","url":"https://api.github.com/u
sers/yhoonkim","avatar_url":"https://avatars.githubusercontent.com/u/6269456?"},"repo":{"id":28594770,"name":"yhoonkim/GraphBoard
 url":"https://api.github.com/repos/yhoonkim/GraphBoard"},"payload":{"action":"opened","issue":{"url":"https://api.github.com,"
                   eact to articles","user":{"login":"yhoonkim","id":6269456,"avatar url":"https://avatars.githubusercontent.com
                 abels":[],"state":"open","locked":false,"assignee":null,"milestone":null,"comments":0,"created_at":"2015-01-01T1
          Join\n\n- [ ] Own board\n\n- [ ] Interview with people who want to archieve own thought within own writings."}},"public
 'id":"2489651096",<mark>"type":"PullRequestEvent"</mark>,"actor":{"id":10357835,"login":"mevlan","gravatar_id":"","url":"https://api.github.c
om/users/mevlan","avatar_url":"https://avatars.githubusercontent.com/u/10357835?"},"repo":{"id":28668460,"name":"mevlan/script","
url":"https://api.github.com/repos/mevlan/script"},"payload":{"action":"opened","number":3,"pull_request":{"url":"https://api.gi
-UUU:**--F1 tmp.json
                             Top L1
```

semi-structured

35+ event types

payload structure and size varies by event type





## **Dataframes and Dynamic Frames**



#### **Dataframes**

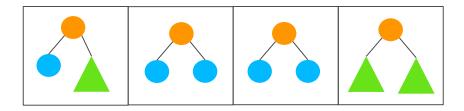
Core data structure for SparkSQL

Like structured tables

Need schema up-front

Each row has same structure

Suited for SQL-like analytics



#### Dynamic Frames

Like dataframes for ETL

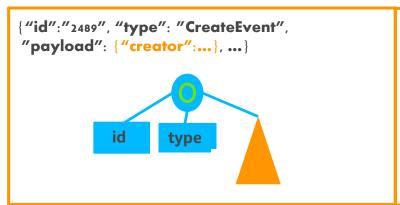
Designed for processing semi-structured data, e.g. JSON, Avro, Apache logs ...

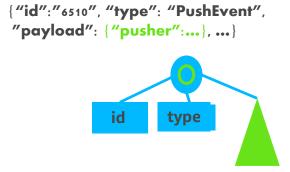


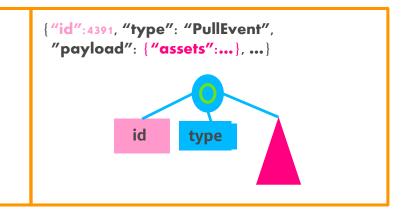


## Dynamic Frame internals

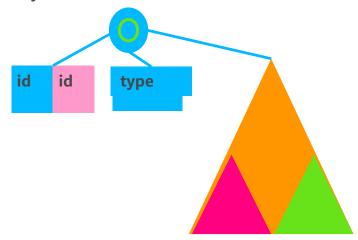
#### Dynamic Records







#### Dynamic Frame Schema



#### Schema per-record, no up-front schema needed

- Easy to restructure, tag, modify
- Can be more compact than dataframe rows
- Many flows can be done in single-pass





## Dynamic Frame transforms

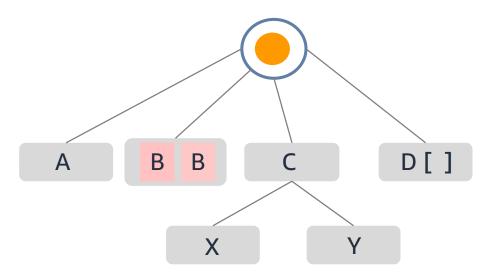
ApplyMapping() A X Y



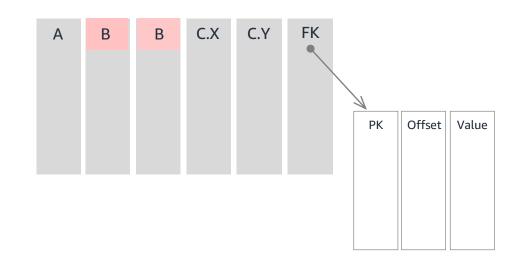


## Relationalize() transform

#### Semi-structured schema



#### **Relational schema**



Transforms and adds new columns, types, and tables on-the-fly

Tracks keys and foreign keys across runs

SQL on the relational schema is orders of magnitude faster than JSON processing





## Useful AWS Glue transforms

toDF(): Convert to a Dataframe

Spigot(): Sample data of any Dynamic Frame to S3

Unbox(): Parse string column as given format into Dynamic Frame

Filter(), Map(): Apply Python UDFs to Dynamic Frames

Join(): Join two Dynamic Frames

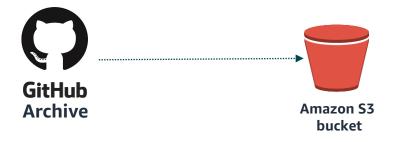
And more ....





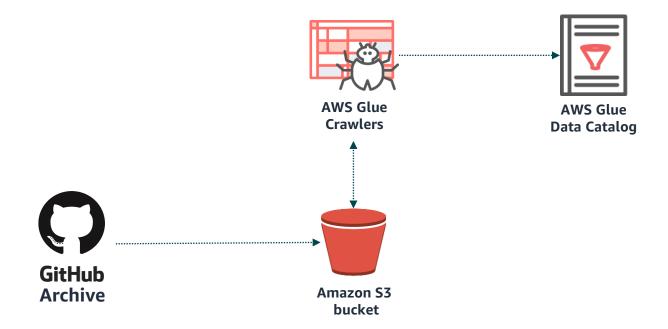






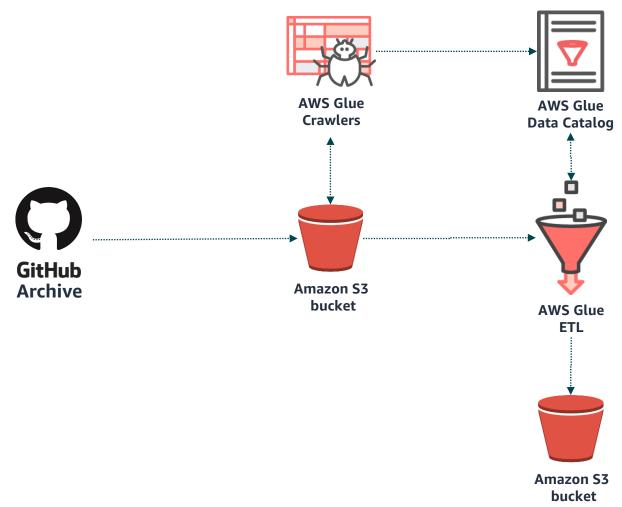






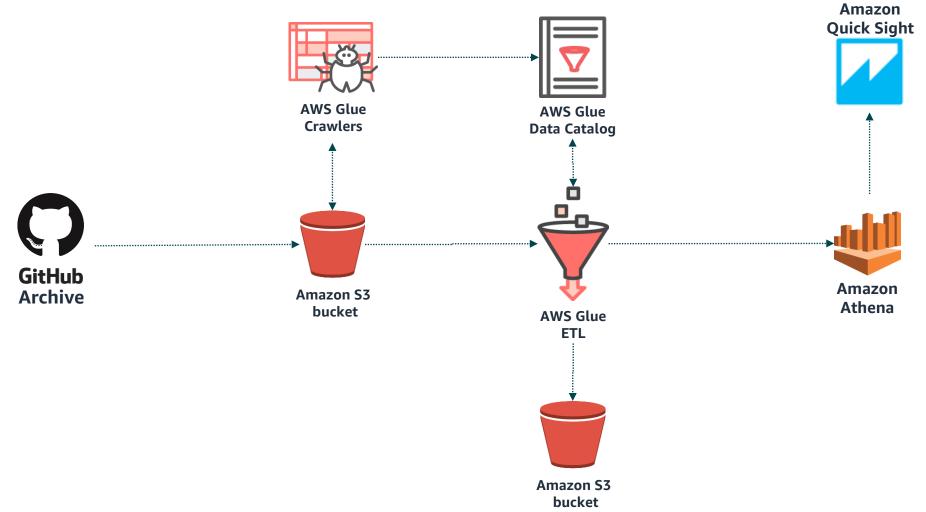


















## Take the demo home...



http://bit.ly/aws-innovate-2018-glue-demo





## NETH Traffic Information Provisioning with AWS Data Lake

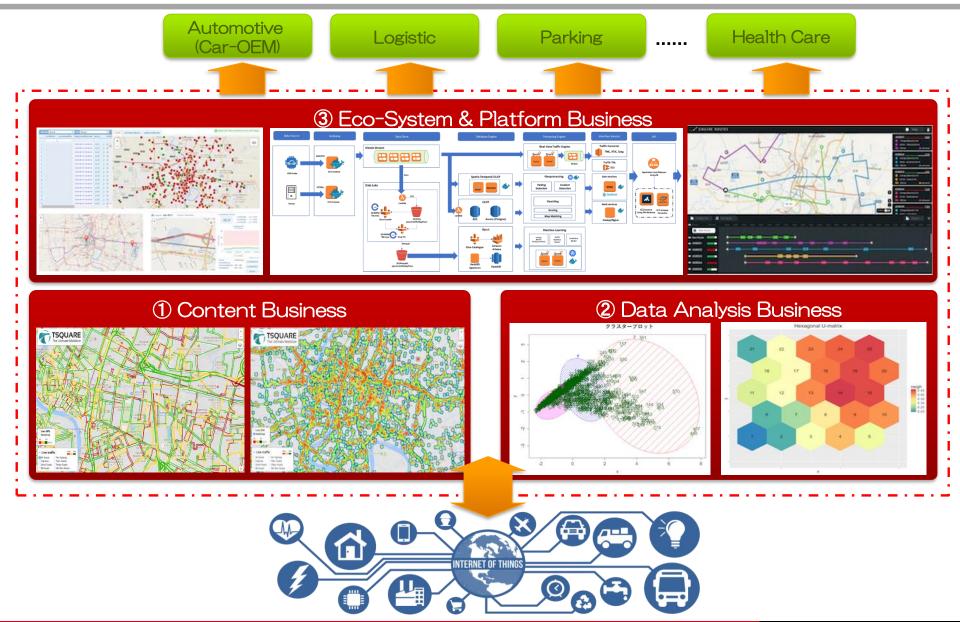


Content Development & Distribution

Thanomsak Ajjanapanya Group Manager Content Department

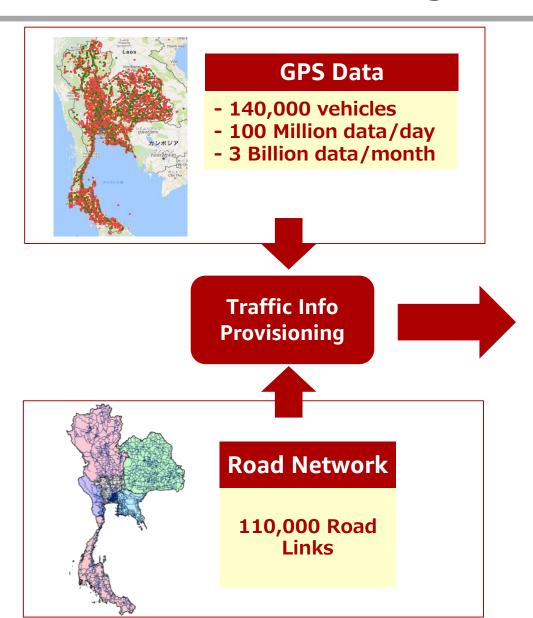
#### **NETH Contents Business Overview**

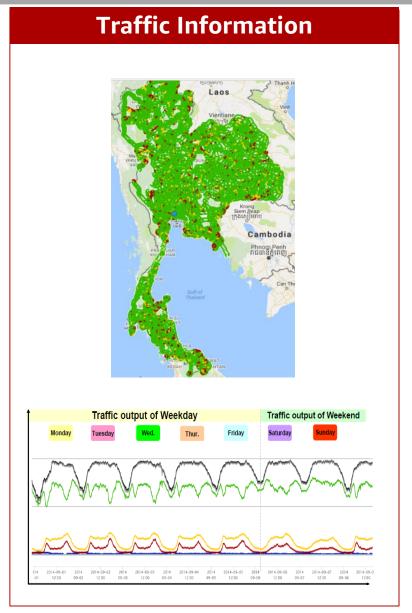




#### **NETH Traffic Provisioning Story**



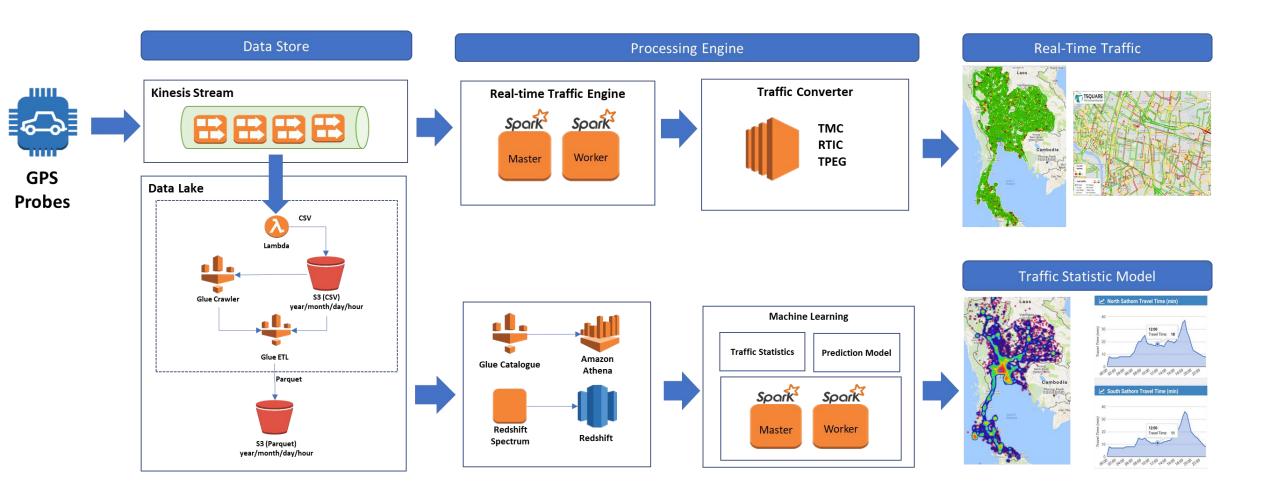




Copyright © NEXTY Electronics Corporation

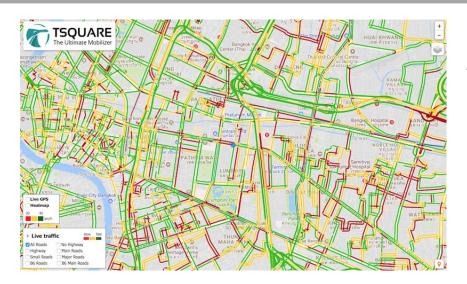
#### **Data Lake Architecture**

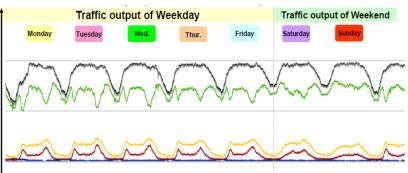




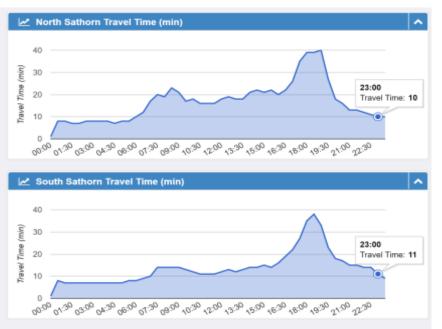
#### **Learning and practices**











## Take the demo home...



http://bit.ly/aws-innovate-2018-glue-demo



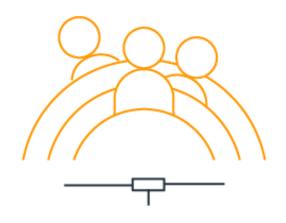


# Learn from AWS experts. Advance your skills and knowledge. Build your future in the AWS Cloud.



#### **Digital Training**

Free, self-paced online courses built by AWS experts



#### **Classroom Training**

Classes taught by accredited AWS instructors



#### **AWS Certification**

Exams to validate expertise with an industry-recognized credential

Ready to begin building your cloud skills?

Get started at: <a href="https://www.aws.training/">https://www.aws.training/</a>





# With deep expertise on AWS, APN Partners can help your organization at any stage of your Cloud Adoption Journey.



#### **AWS Managed Service Providers**

APN Consulting Partners who are skilled at cloud infrastructure and application migration, and offer proactive management of their customer's environment.



#### **AWS Competency Partners**

APN Partners who have demonstrated technical proficiency and proven customer success in specialized solution areas.



#### **AWS Marketplace**

A digital catalog with thousands of software listings from independent software vendors that make it easy to find, test, buy, and deploy software that runs on AWS.



#### **AWS Service Delivery Partners**

APN Partners with a track record of delivering specific AWS services to customers.

Ready to get started with an APN Partner?

Find a partner: <a href="https://aws.amazon.com/partners/find/">https://aws.amazon.com/partners/find/</a>

Learn more at the AWS Partner Network Booth





## **Thank You for Attending AWS Innovate**

We hope you found it interesting! A kind reminder to **complete the survey.** 

Let us know what you thought of today's event and how we can improve the event experience for you in the future.

- aws-apac-marketing@amazon.com
- twitter.com/AWSCloud
- facebook.com/AmazonWebServices
- youtube.com/user/AmazonWebServices
- slideshare.net/AmazonWebServices
- twitch.tv/aws



