

A faint, light gray wireframe globe is centered in the background, showing latitude and longitude lines.

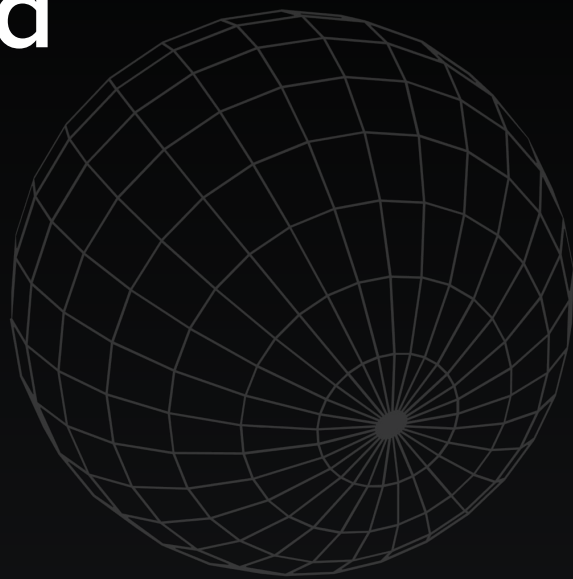
May 23—27, 2022

AIRFLOW SUMMIT

Large, bright green abstract shapes are positioned on the left and right sides of the image. On the left, there is a thick, curved line and a solid rectangular block below it. On the right, there is a thick, curved line. Small green triangles point towards the center from the left and right sides.

Orchestrating hybrid workflows with Apache Airflow

Ricardo Sueiras
Principal Developer Advocate, AWS





Apache
Airflow



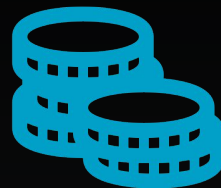
Regulation &
Compliance



Heritage
systems



Reduce
complexity



Cost
effective

SqlToS3Operator

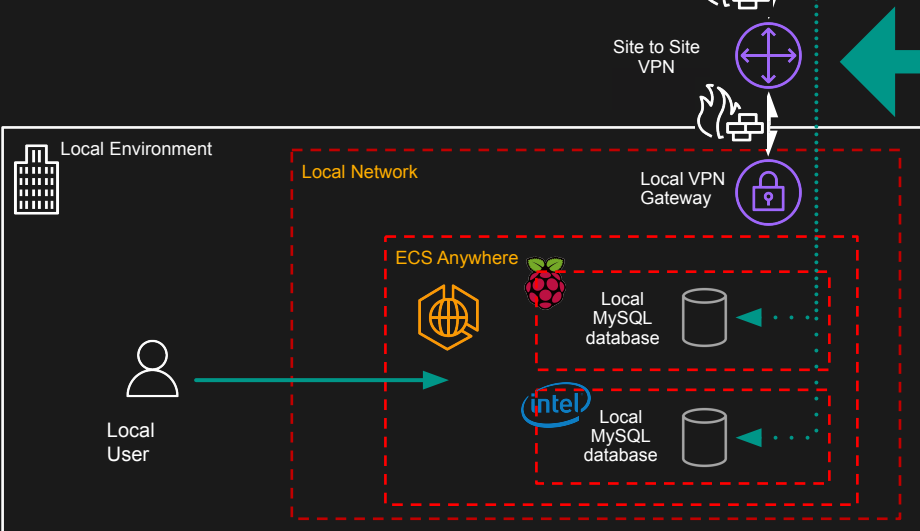
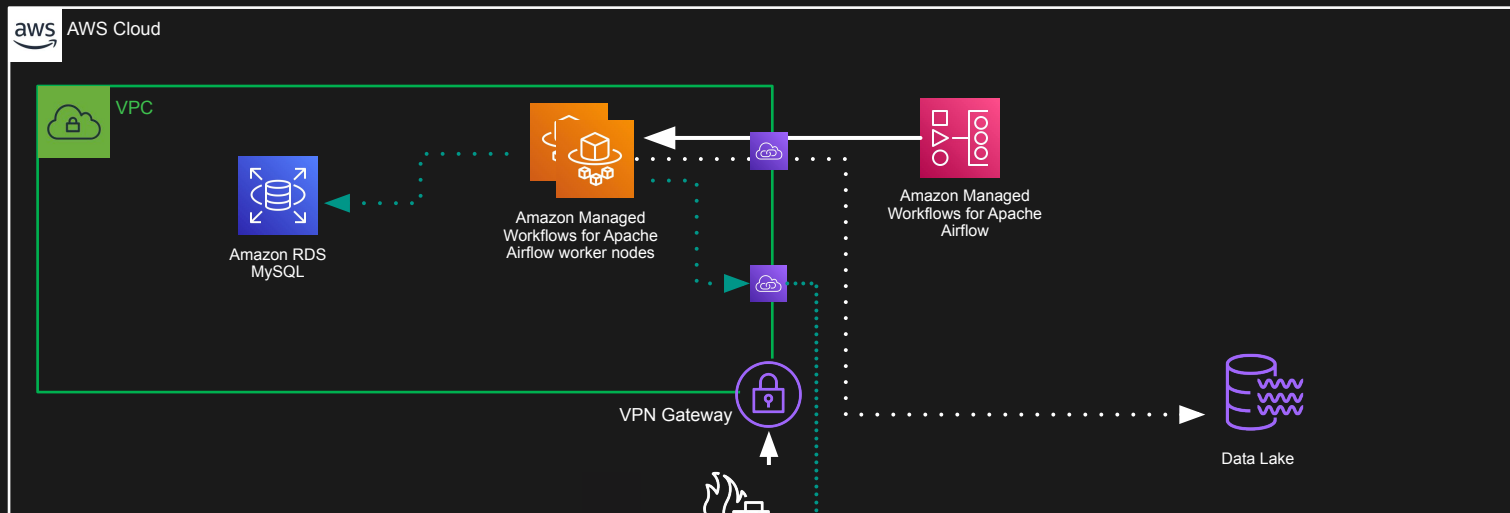
- Need to integrate networks between Cloud and remote (via VPN for example)
- Additional work required to enable connectivity (inbound/controlled networks)
- Processing of data in Airflow Worker

PythonOperator

- Need to integrate networks between Cloud and remote (via VPN for example)
- More complex DAGs
- Processing of data in Airflow Worker
- Potential Apache Airflow anti-pattern

AthenaOperator

- Need to integrate networks between Cloud and remote (via VPN for example)
- Processing of data in Cloud (Athena)
- Create Athena Federated Queries (running as Lambda functions)
- Additional work required to enable connectivity (inbound/controlled networks)

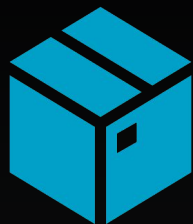


Bridging this gap.

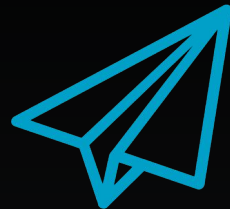
Using existing operators
may require you to deploy
a VPN solution to access
remote locations



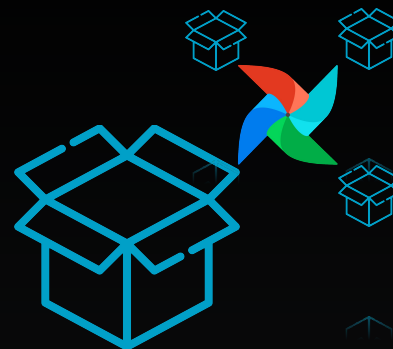
Develop ETL
application



Package as
container



Push to Container
Repository



Orchestrate
running
container

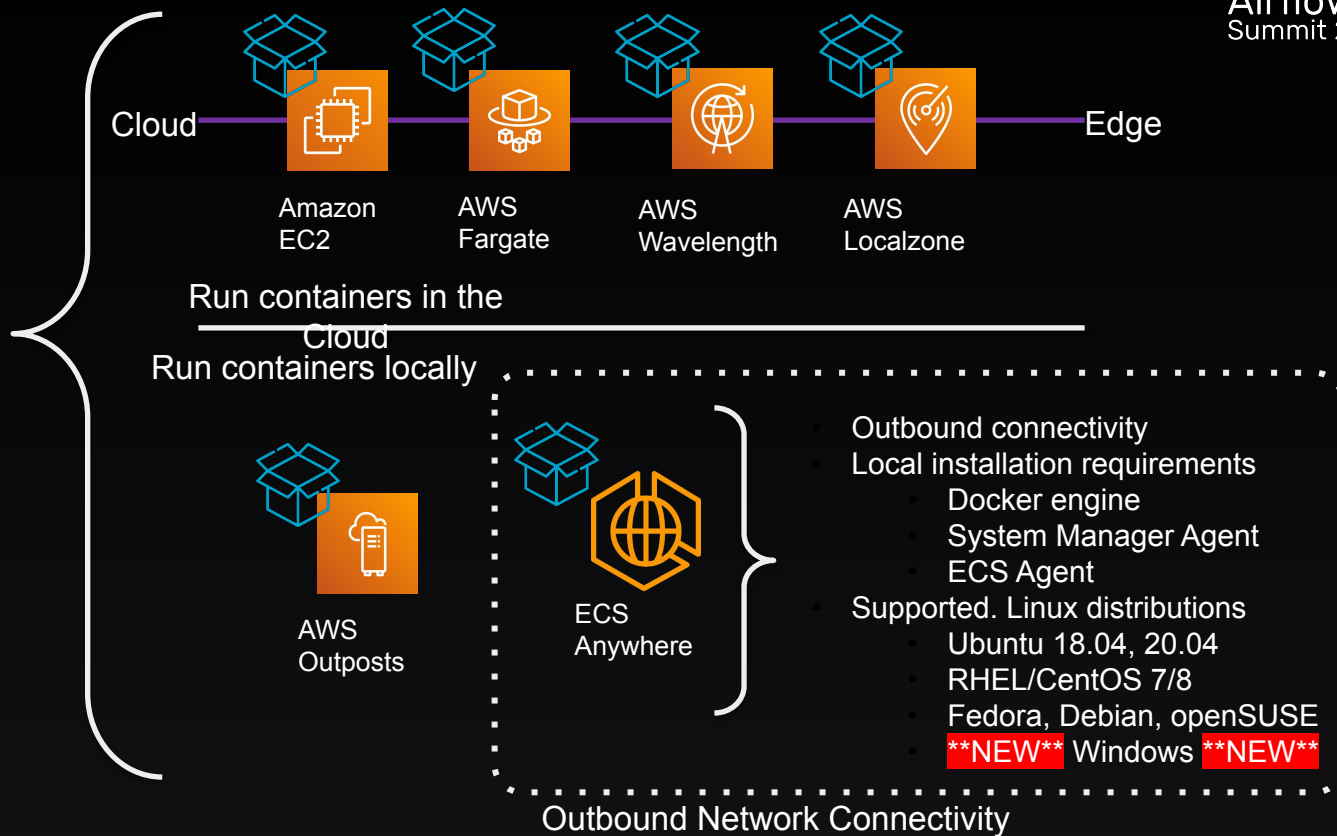
KubernetesOperator

- Need to integrate networks between Cloud and remote (via VPN for example)
- Need to build your ETL container image
- Good re-usability
- Need to provision and manage K8s clusters locally

ECSSOperator

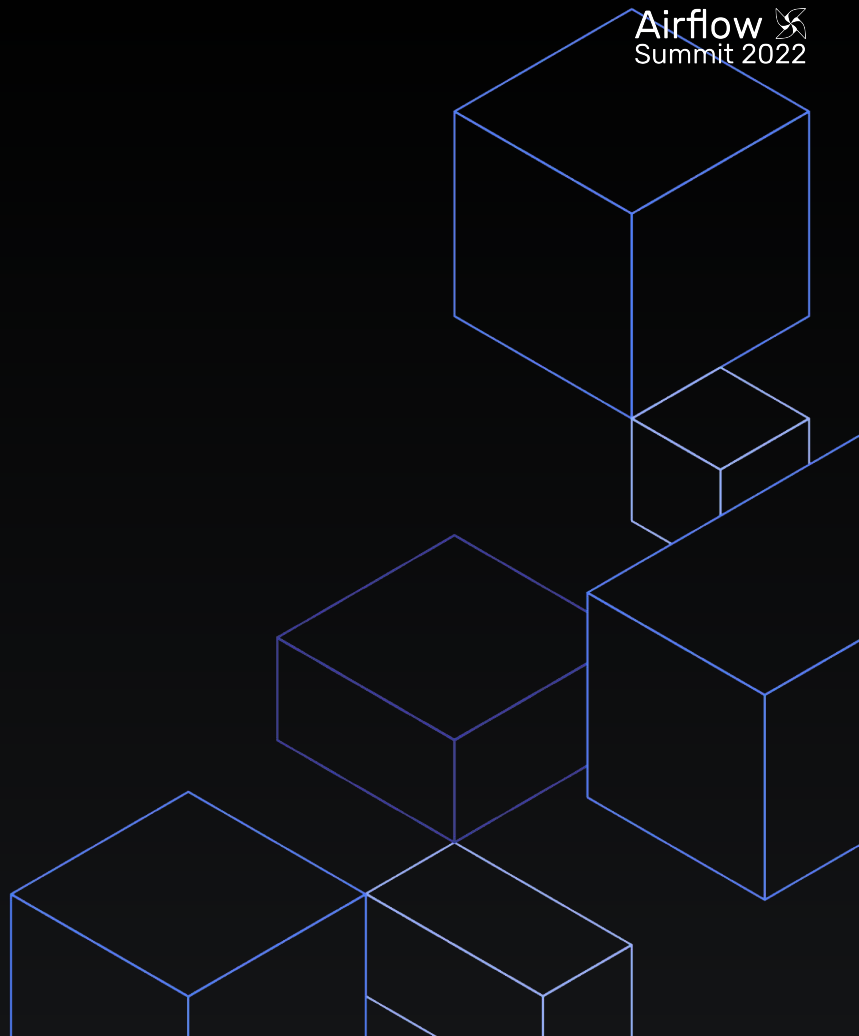
- Need to install local ECS Anywhere agent and create Amazon ECS clusters
- Not open source
- Need to build your ETL container image
- Good re-usability
- Processing of data locally

Run container on ECS Cluster

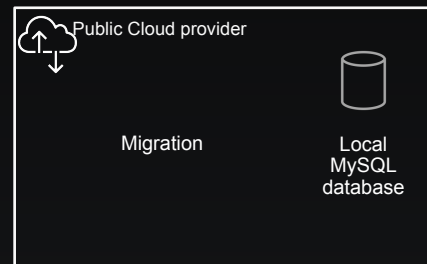
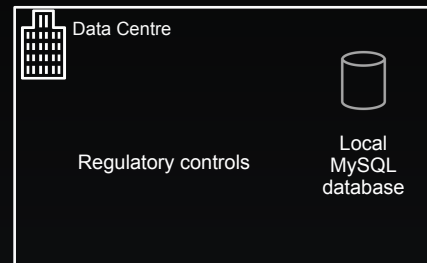
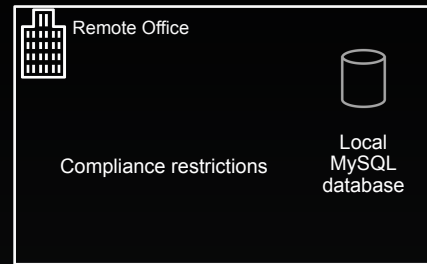
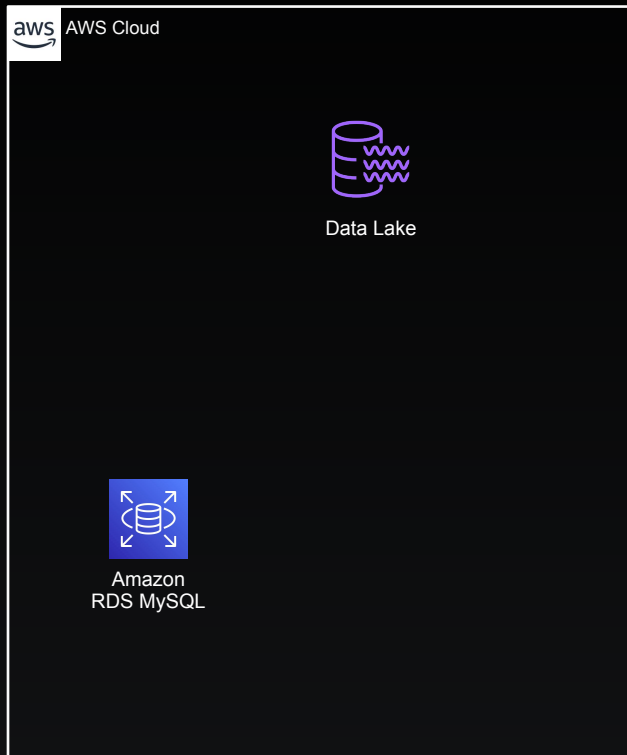


Demo

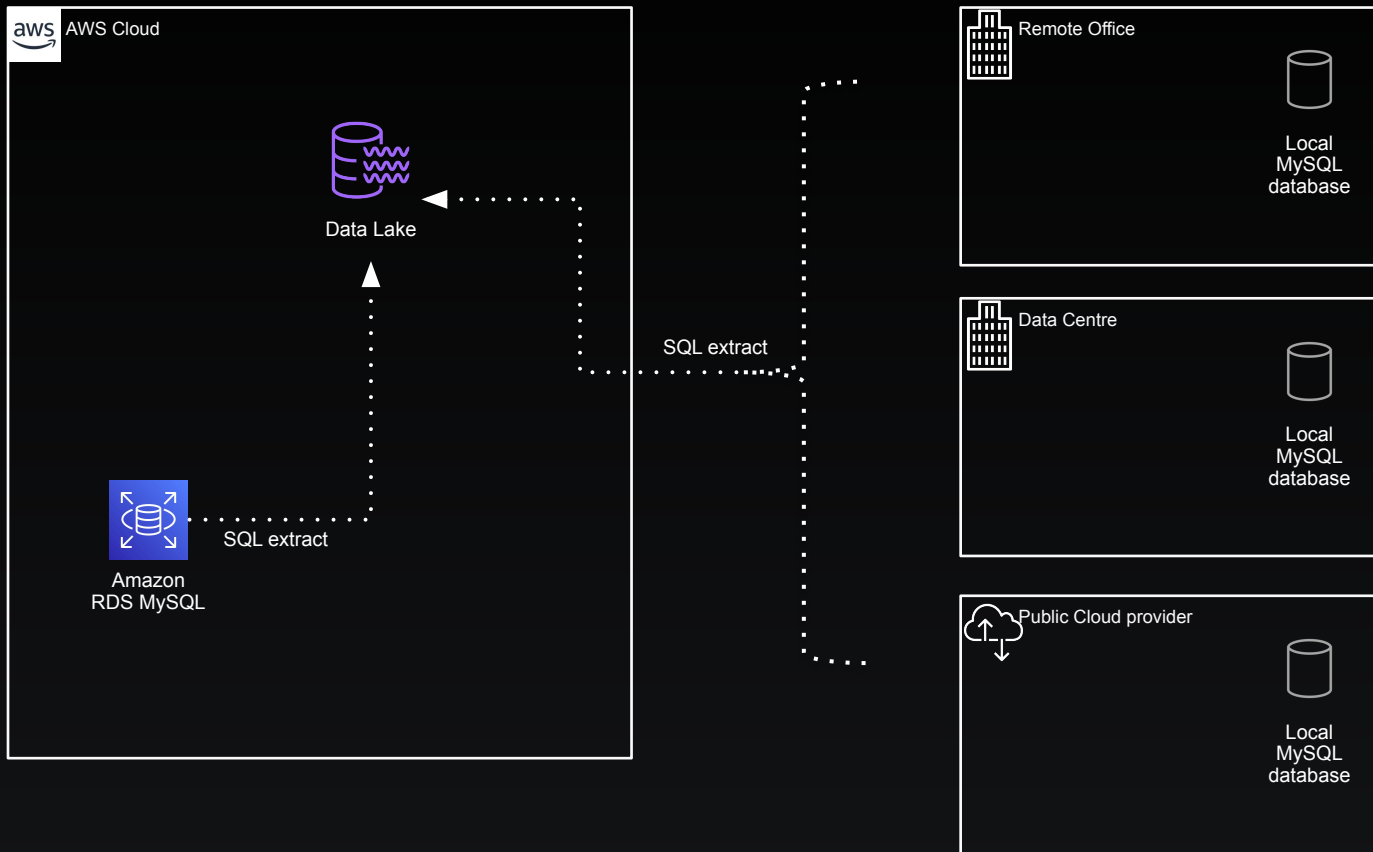
(live, pray for me!)



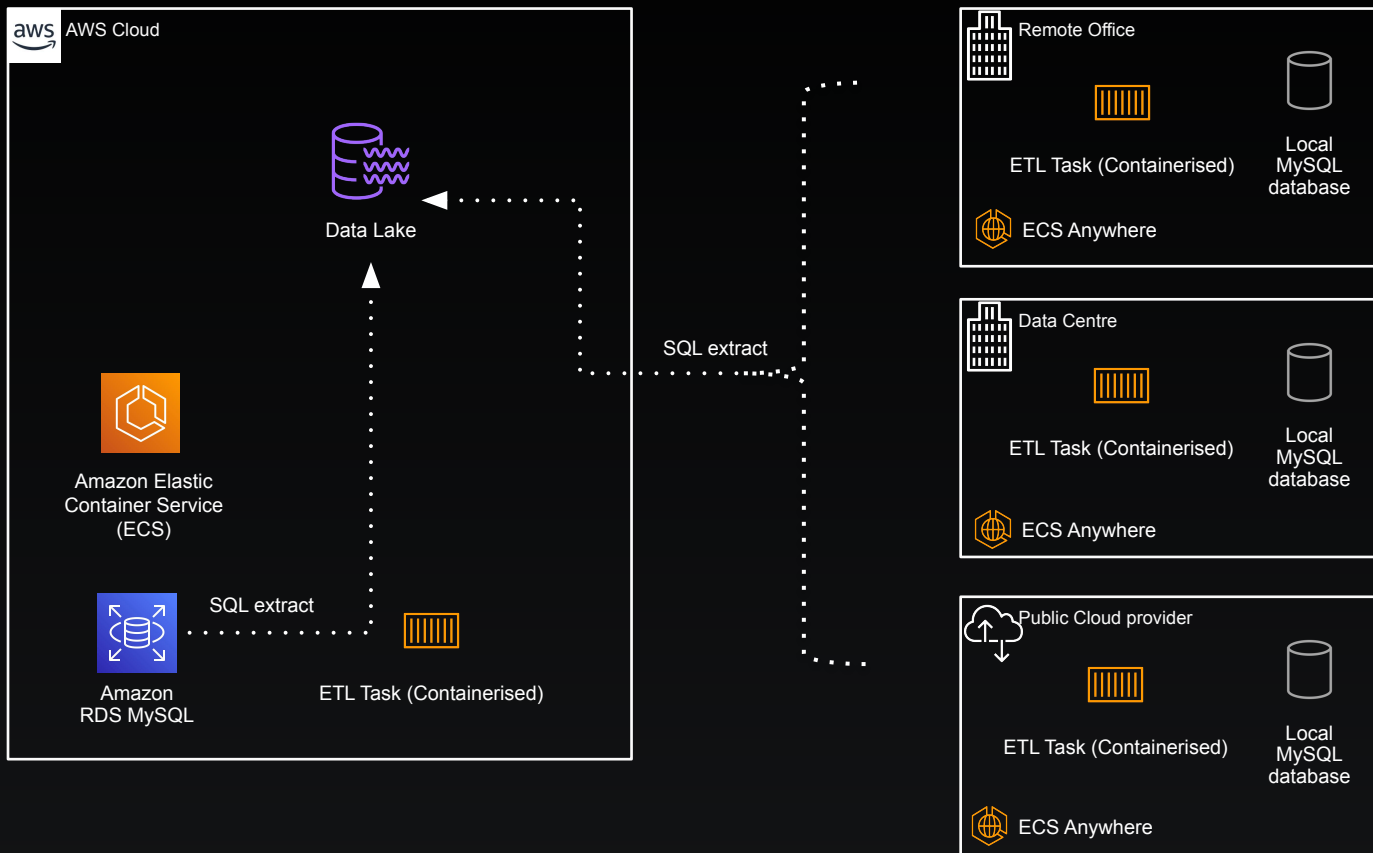
Hybrid data pipeline



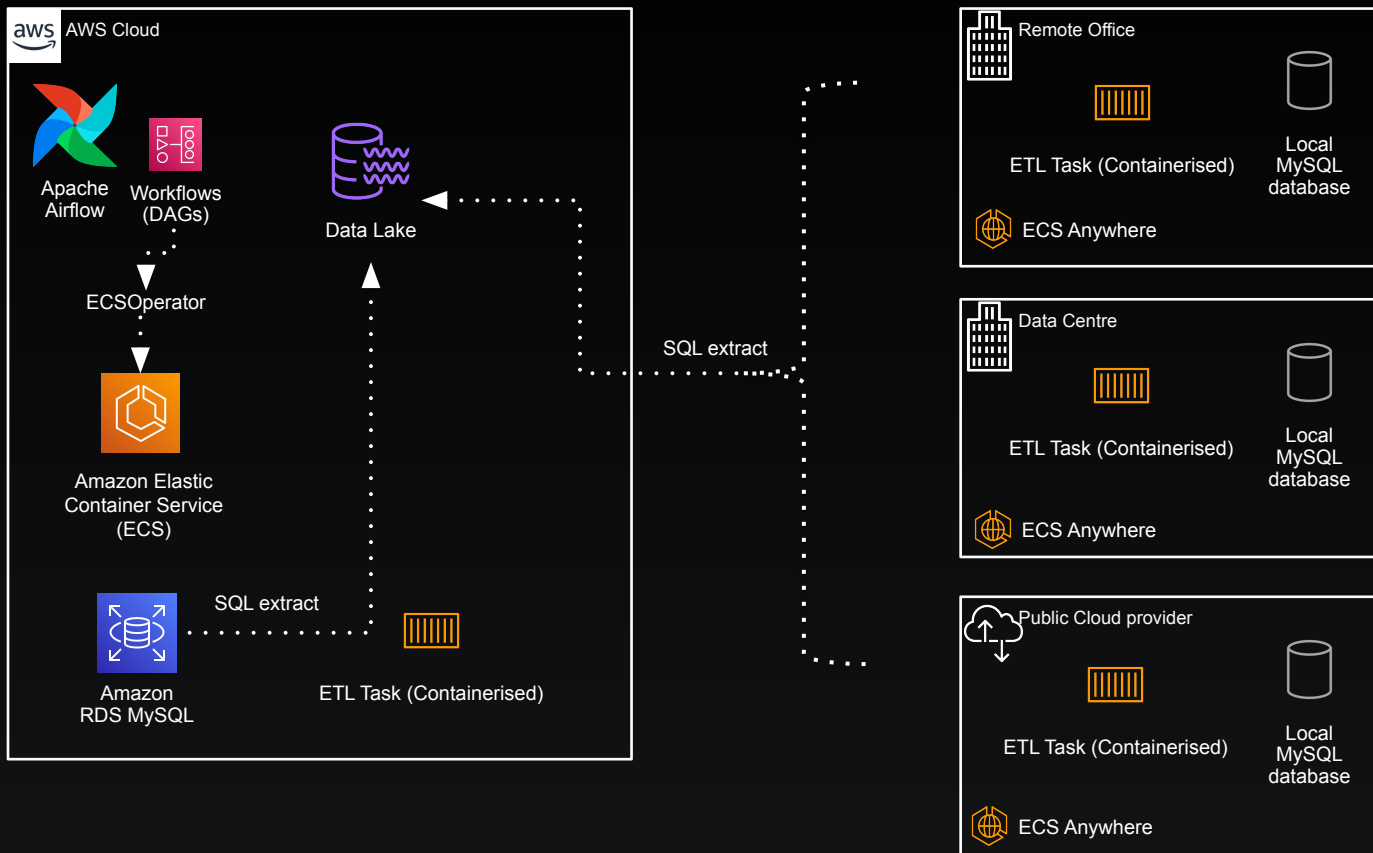
Hybrid data pipeline



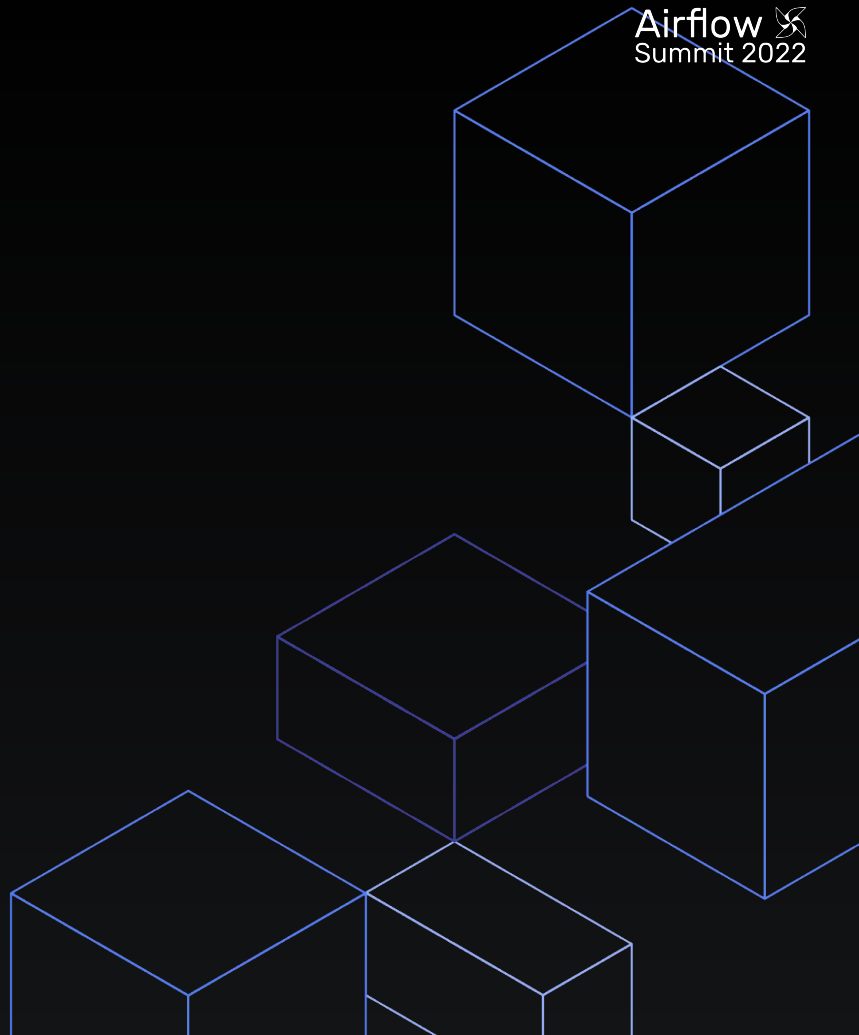
Hybrid data pipeline

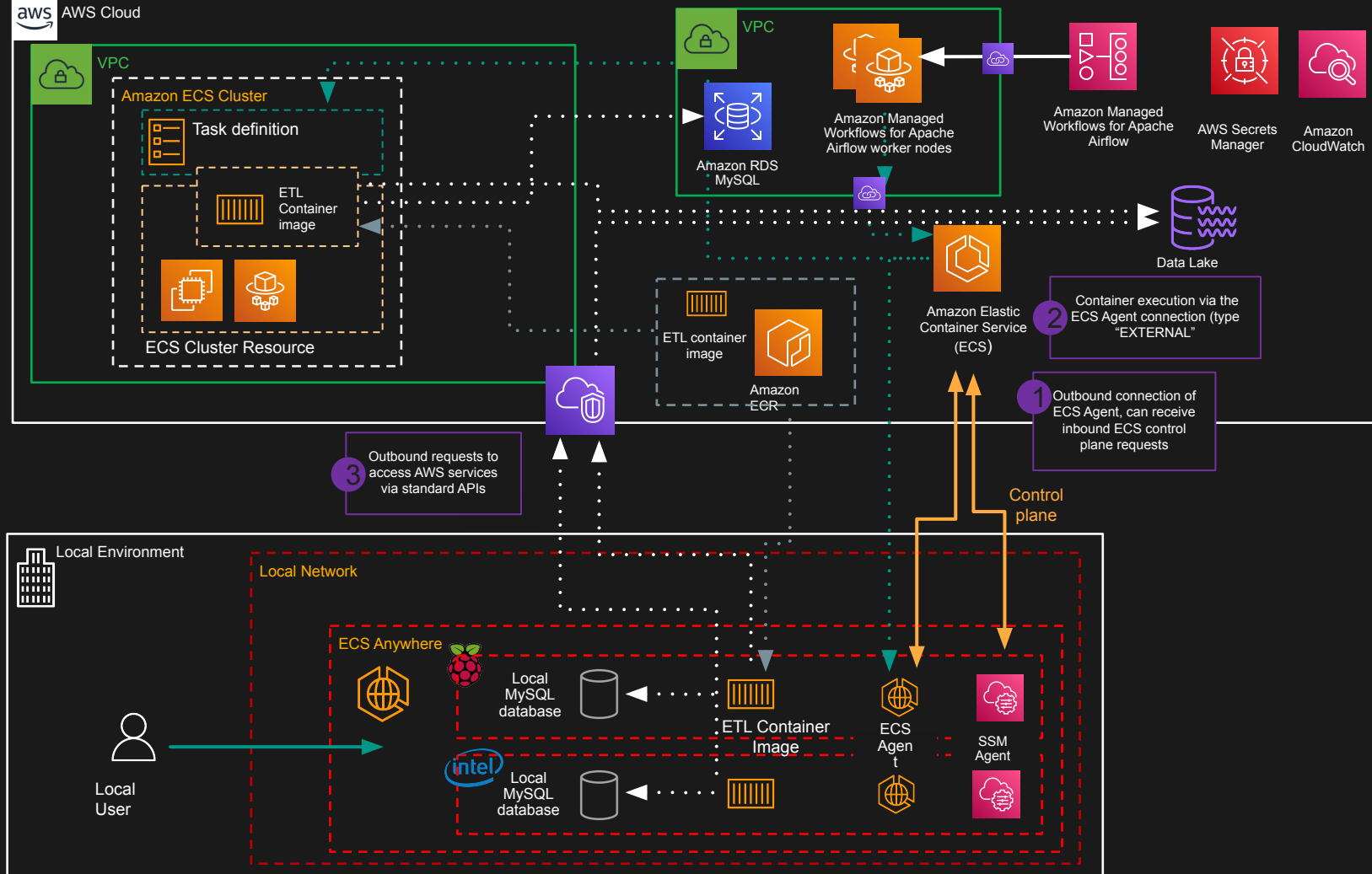


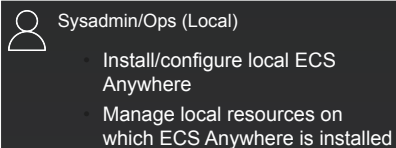
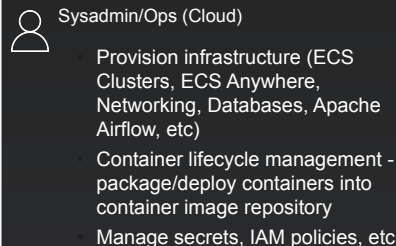
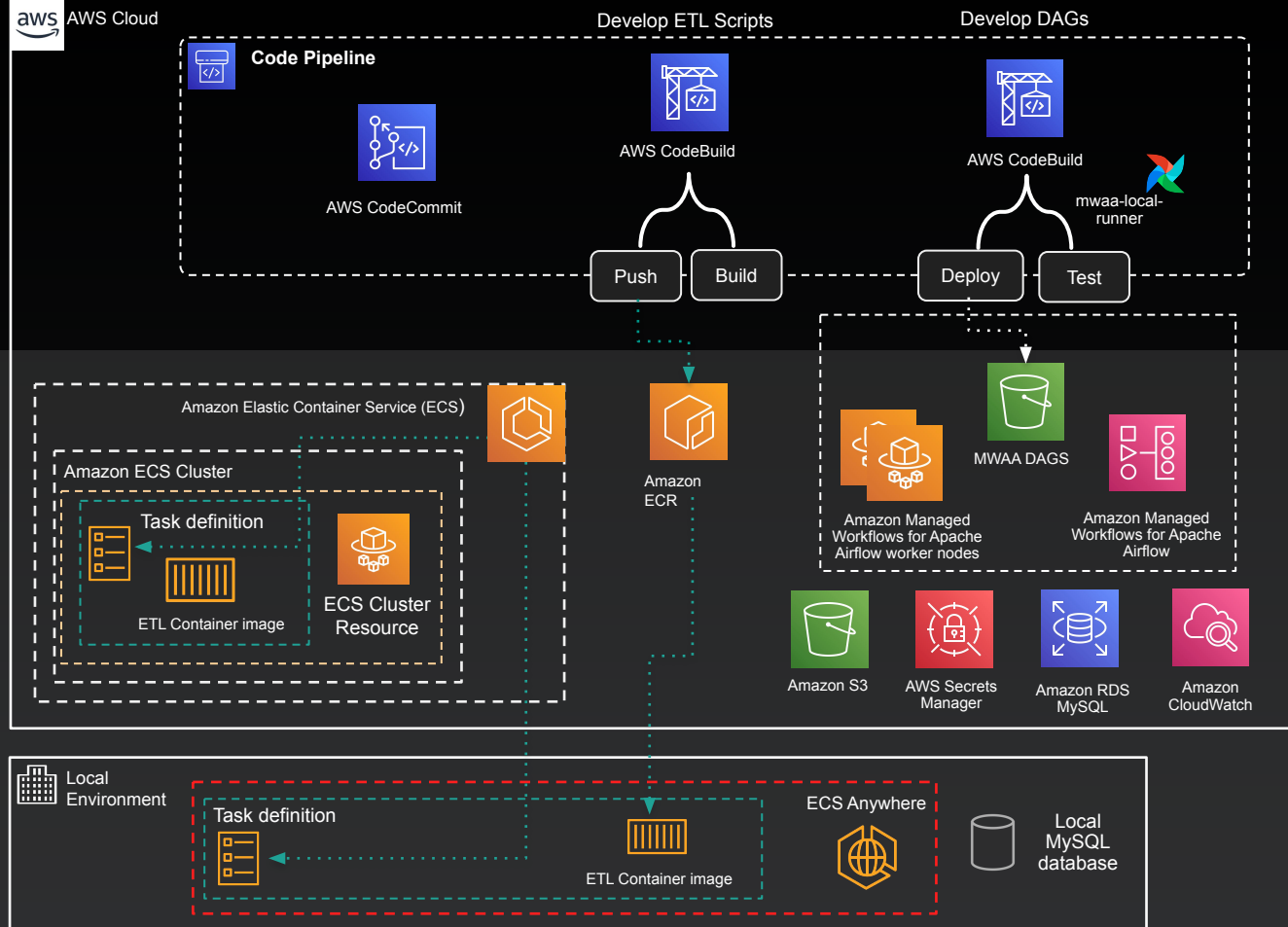
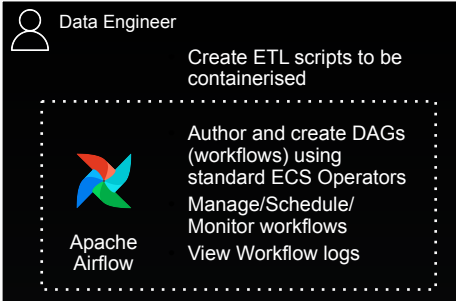
Hybrid data pipeline



Demo







Permissions

Task Definition Role
(`ecsTaskExecutionRole`)

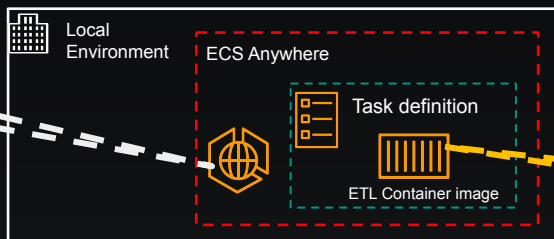
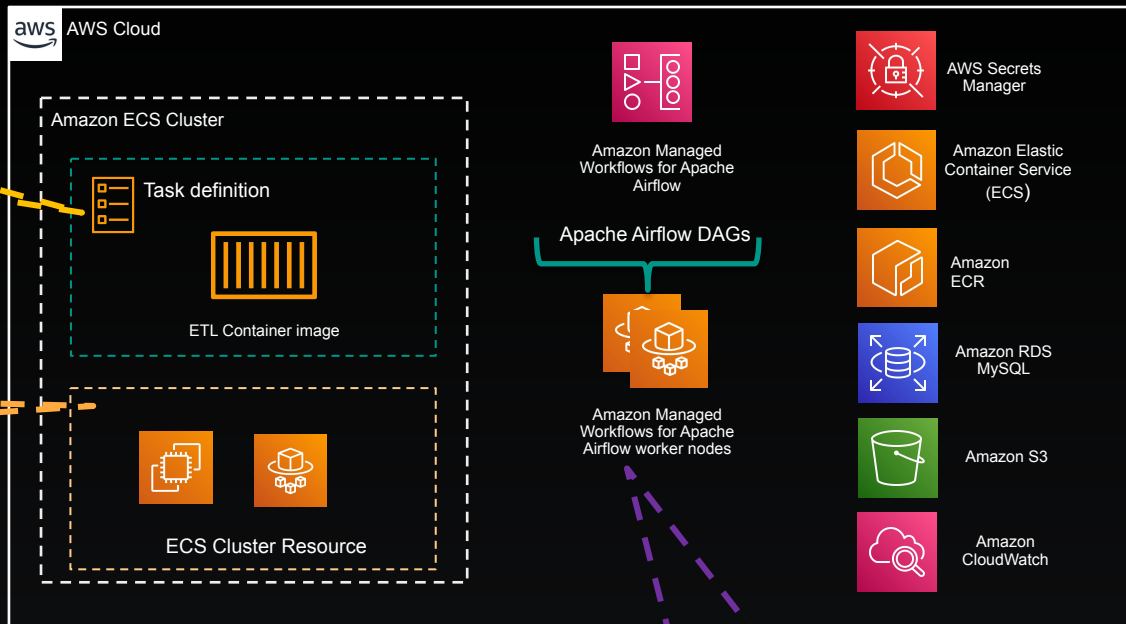
IAM policies needed for your application

Task Execution Role
(`ecsInstanceRole`)

IAM policies needed for the hosts to run your Containers

ECS Anywhere Task Execution Role

Like the Task Execution Role but for the ECS Anywhere agent



MWAA Execution Role

IAM policies needed to access AWS services from your DAGs

Task Definition Role
(`ecsTaskExecutionRole`)

IAM policies needed for your application



<https://github.com/094459/blogpost-airflow-hybrid>

Thank you!

Ricardo Sueiras



<https://www.linkedin.com/in/ricardosueiras>



@094459