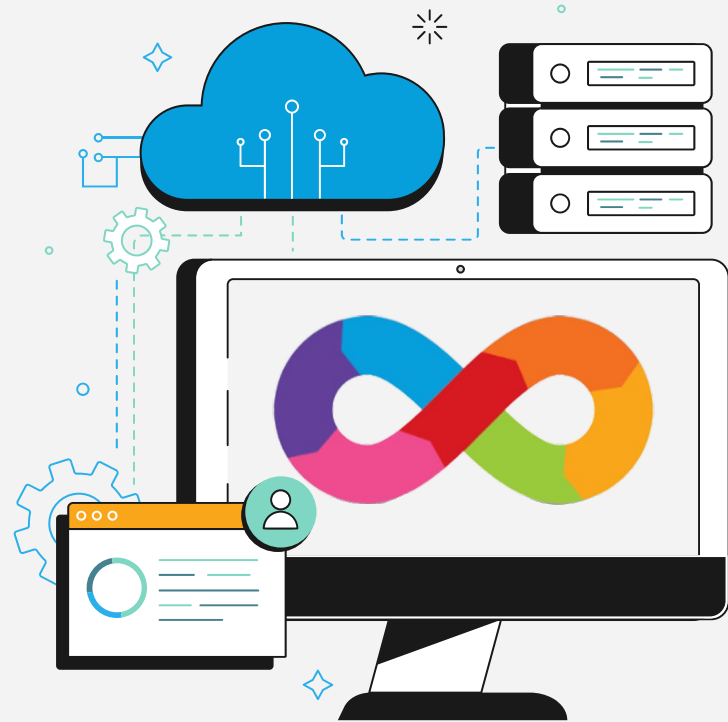


Introduction to DevOps

@ IBA – SMCS

Week 14 – 2
**Disaster Recovery &
MLOPs**



Obaid ur Rehman
Software Architect / Engineering Manager @ Folio3

What is a Disaster Recovery

- Any event that has negative impact on a business continuity or finance is a Disaster.
- Disaster Recovery (DR) is all about recovering from a disaster and resume normal operations.
- Let's discuss two terms: RPO & RTO

RPO & RTO

RTO

The recovery time objective (RTO) is the targeted duration of time between the event of failure and the point where operations resume.

RPO

A recovery point objective (RPO) is the maximum length of time permitted that data can be restored from, which may or may not mean data loss.

RPO & RTO

Business continuity

How much data can you afford to recreate or lose?

**How quickly must you recover?
What is the cost of downtime?**



Disaster Recovery Strategies

- Backup and Restore
- Pilot Light
- Warm Standby
- Hot site / Multi site approach

RTO & RPO of Disaster Recovery Strategies

Backup & Restore

RTO/RPO:
Hours

- Lower priority use cases
- Solutions: Cloud Storage, Backup Solutions
- Cost: \$ to \$\$

Pilot Light

RTO/RPO:
10s of Minutes

- Lower RTO/RPO requirements
- Solutions: Database Service, Replication Solutions
- Cost: \$\$

Warm Standby

RTO/RPO:
Minutes

- Core Applications and Services
- Solutions: Cloud Storage, Database Service, Replication Solutions
- Cost: \$\$\$

Multi Site

RTO/RPO:
Real-Time

- Mission Critical Applications and Services
- Solutions: Database Service, Replication Solutions
- Cost: \$\$\$\$

Backup & Restore

- Cheapest and oldest method of DR Strategy.
- High RPO & possibly High PTO.
- Simple case could be: You take backup of database. All infrastructure is IaC. Recovery would be a manual process.

Pilot Light

- In this approach, we basically keep a minimal set of resources running in another environment or region, which will be ready to go live. We can spin it up when it's needed.
- Whatever takes longest to spin up is kept in stand by mode and rest of the infrastructure we spin when needed. For example, database takes the longest to spin up in another env or to restore its backup, then that needs to be kept up with a small machine but in sync with the prod data, so, when needed, we scale the system up and it's ready to go live. Rest of the infrastructure can be spin up using IaC.

Pilot Light

- It's cheap to run this minimal infrastructure (+operational costs) and we are not paying for the rest of the infrastructure, until we need them.

Warm Standby

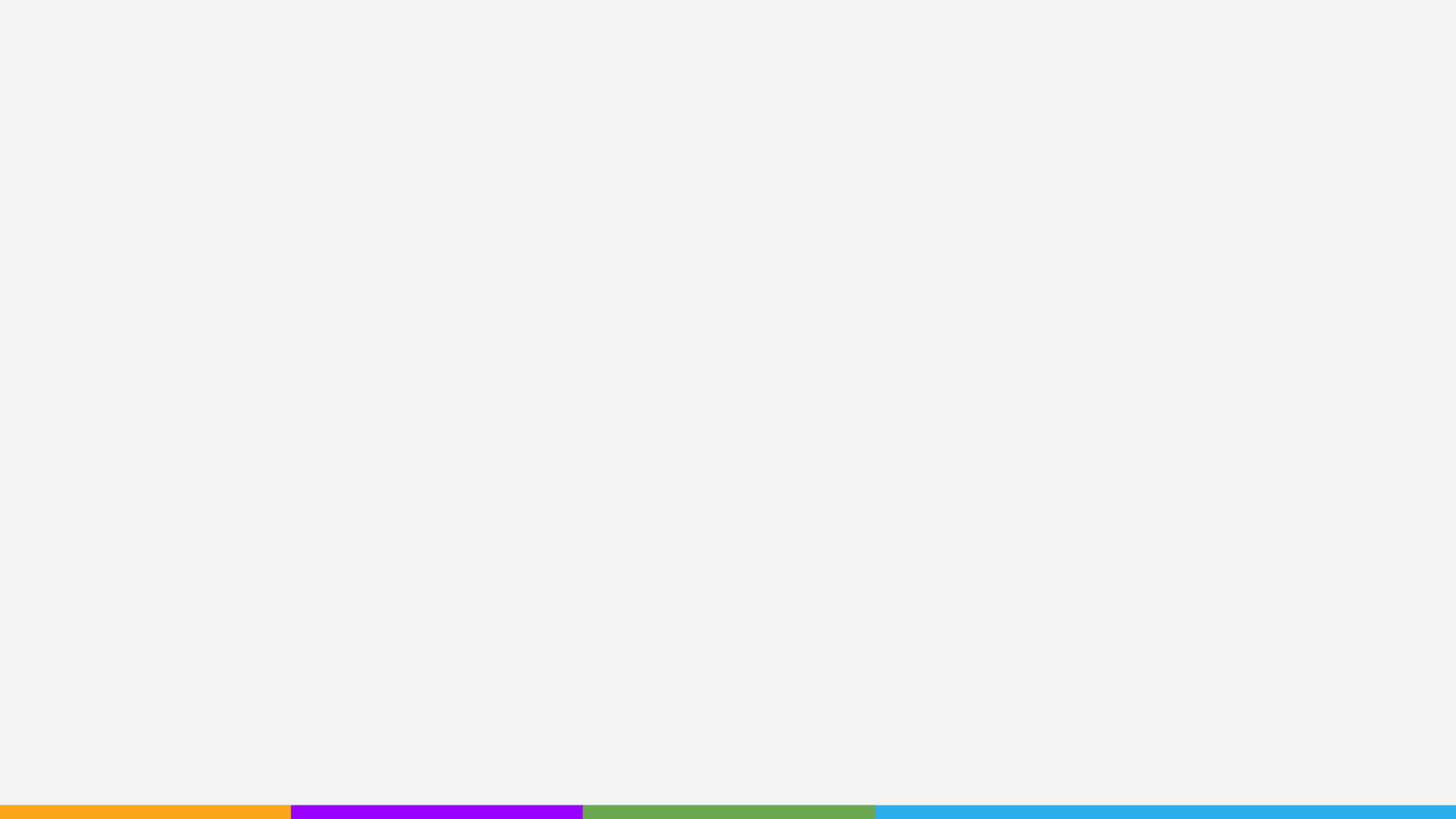
- Kind of similar to Pilot light but additional services are duplicated.
- Full System is duplicated but at a minimum size.
- Upon Disaster we scale it up.

Multi-site (active – active)

- This is the costliest approach of all the solutions, but have the shortest RTO/RPO.
- We basically have identical copies of production infrastructure which is running 24 * 7 side by side both in active status fronting by a load balancer. Best suitable for business generating critical applications.

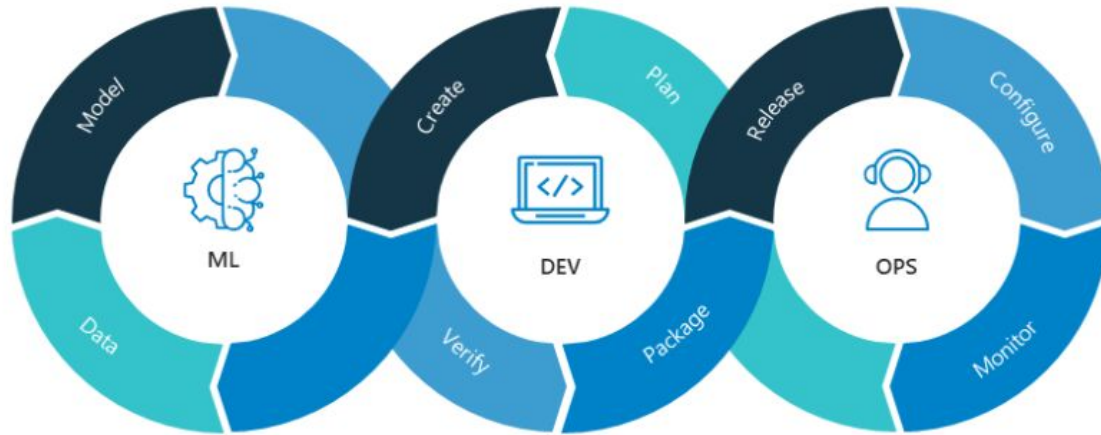
Conclusion

- Whatever the solution we choose, that needs to be tested thoroughly to ensure, that, during actual DR situation, it works.
- Monitoring tools tied up with alerts/events, using which, remediate the application back to its desired state automatically, is a very significant/elegant process, which requires a lot of testing. Automating these strategies saves a lot of time, which directly means meeting RTO/RPO SLAs.



What is MLOps

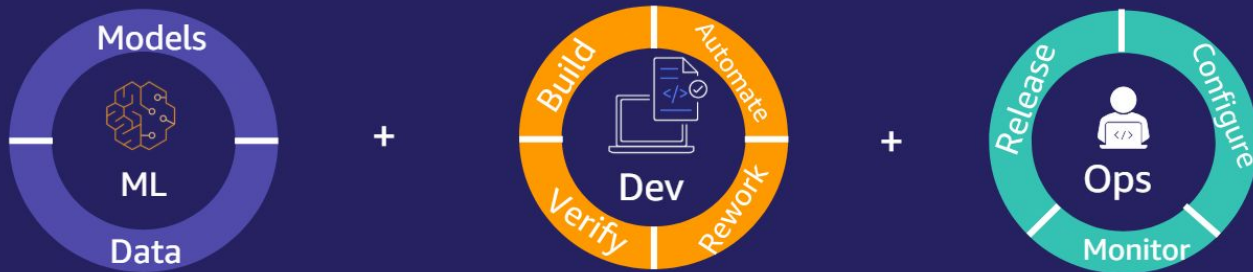
Machine learning operations (MLOps) are a set of practices that automate and simplify machine learning (ML) workflows and deployments.



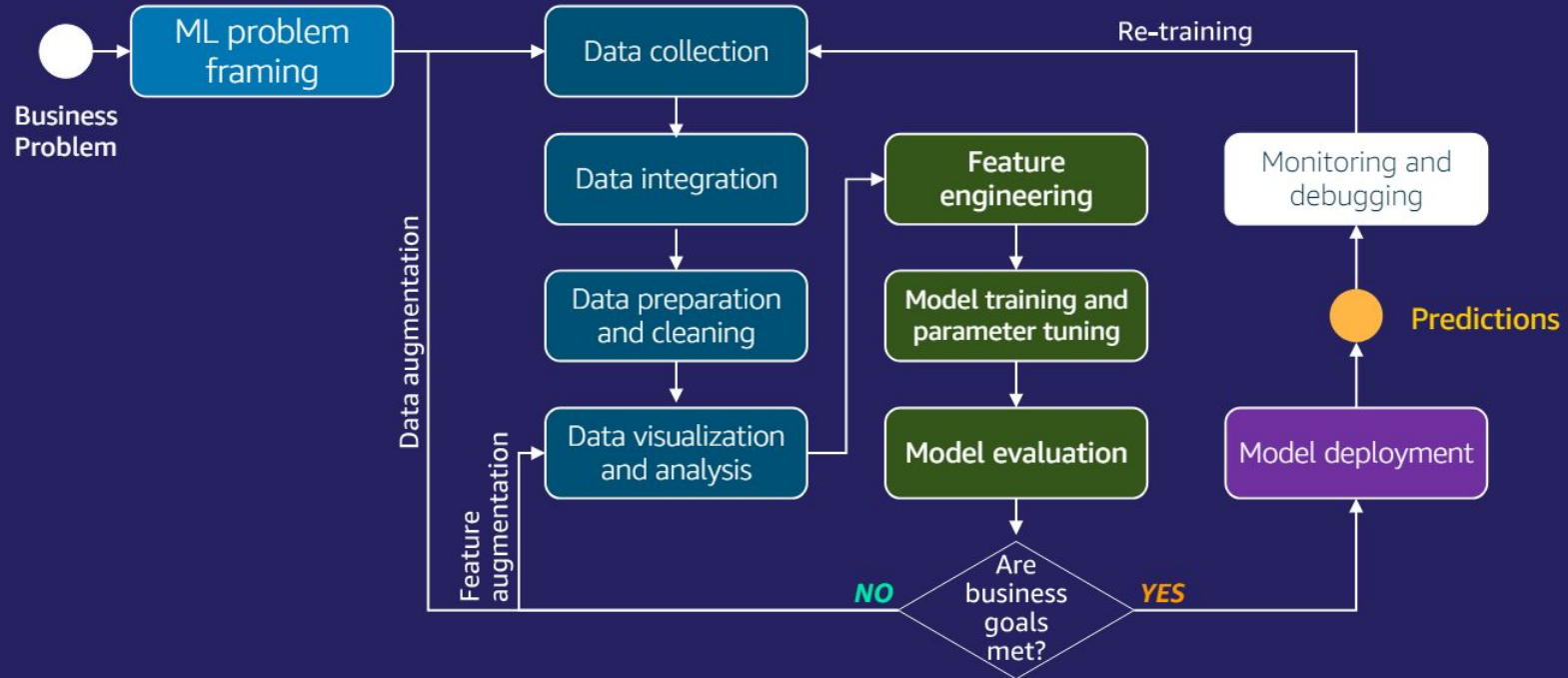
What is MLOps

ML + Dev + Ops = MLOps

Collaborative and experimental in nature | Automate as much as possible |
Continuous improvement of ML Models | Standardize and Scale

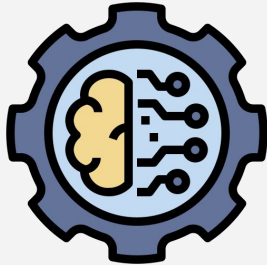


The ML Process



The Friction

- Build
- Train
- Deploy
- Monitor
- Manage
- Re-train

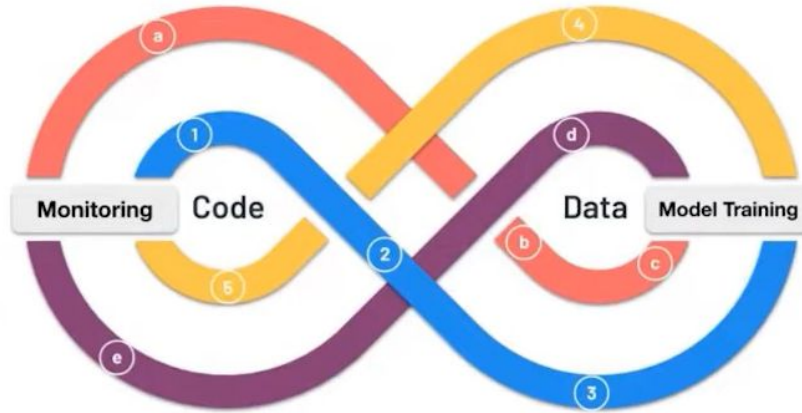


Manual Process can become bottleneck, impact productivity and become costly.

Only 53% of POCs make it into production
(Source: Gartner)

75% of organizations will shift from piloting to operationalizing AI
(Source: Gartner)

Code and Data Loops



Code Loop:

Iterations on machine learning code

Dev

- 1 Define opportunity and desired outcomes
- 2 Design and prototype ML techniques
- 3 Model training, tuning, and selection

Prod

- 4 Testing, integration, and productionization
- 5 Deployment, system and prediction monitoring

Data Loop:

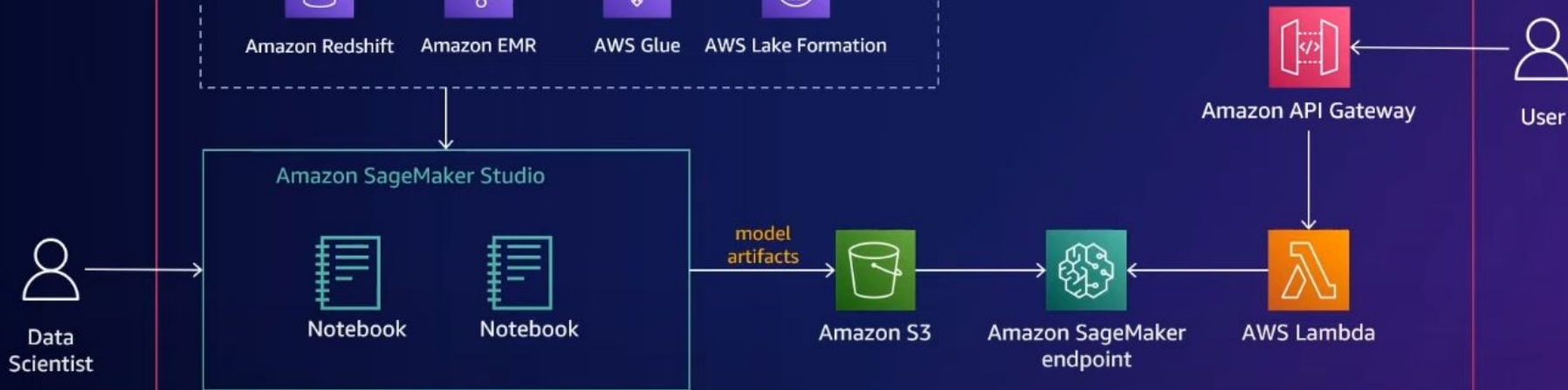
Developments on data

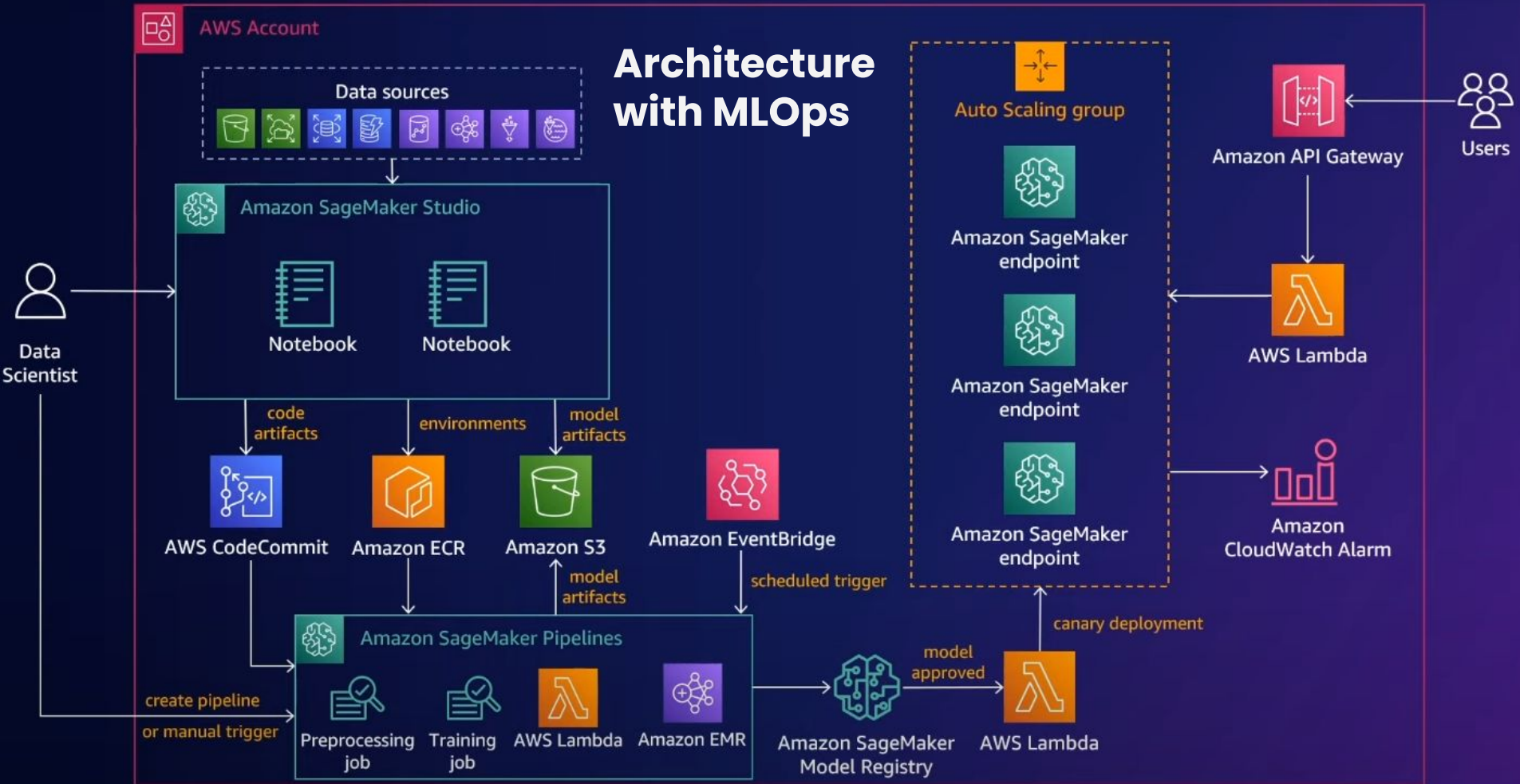
- a Identify and collect data from sources
- b Clean and label data
- c Explore, prepare, and split data
- d Evaluate test data performance
- e Training vs. production data monitoring



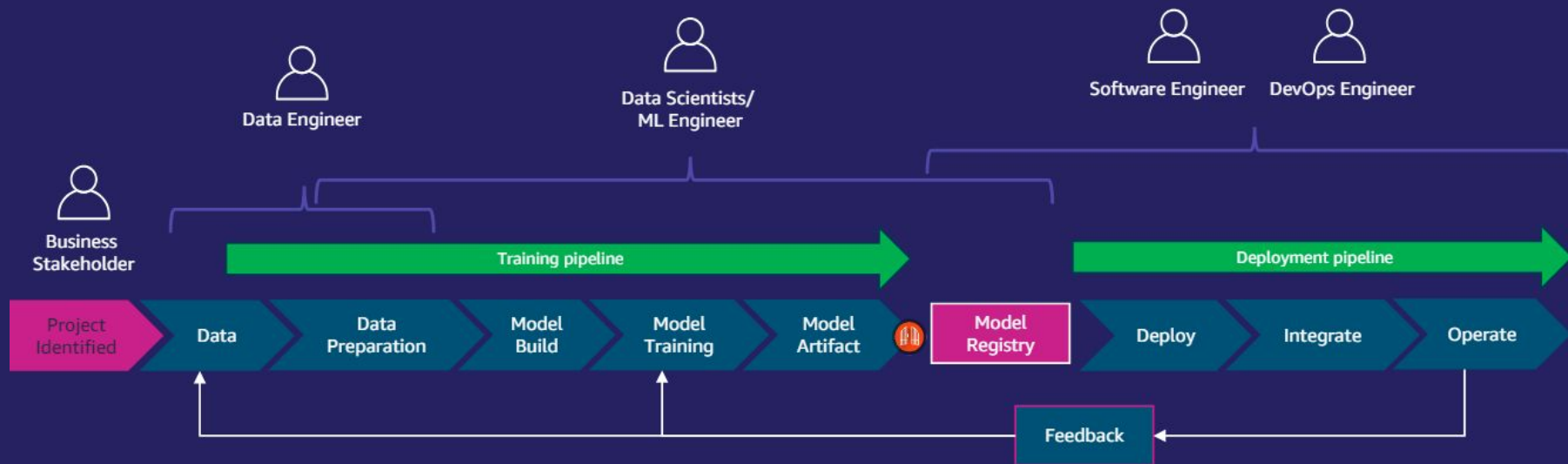
AWS Account

Architecture without MLOps





ML Ops Practices



DevOps V/s MLOps

	DevOps	MLOps
Main Purpose	Automation of software processes like quality assurance and feedback loops	Standardize Machine Learning Lifecycle by processes and automated quality checks
Deployment cycles	Frequent stepwise iterations	Continuous, long training cycles
Team composition	Software Development + DevOps + QA Team	Machine Learning Engineers + Data Engineers/Scientist
Deliverable	Code and integration	Model + Training Data + Training parameters
Objectives	Business goals	Exploration of data and Model Experiments

Whats next: LLMOPs

LLMOPS: Large Language Model Ops encompasses the practices, techniques and tools used for the operational management of large language models in production environments.



End

Q&A

