



Building Your Data Mesh with SageMaker Unified Studio

A Practical Implementation Blueprint for Multi-Account,
Automated Environments

A tailored guide for implementing a robust, CI/CD-driven data mesh architecture using on-premises Git and a custom governance application.

The Objective: An Automated, Governed Data Mesh Ecosystem

We are tasked with building a SageMaker Unified Studio (SMUS) environment that supports a true data mesh operating model. This requires a solution that is not just functional, but automated, secure, and scalable across distinct development lifecycles.

Core Requirements Checklist:

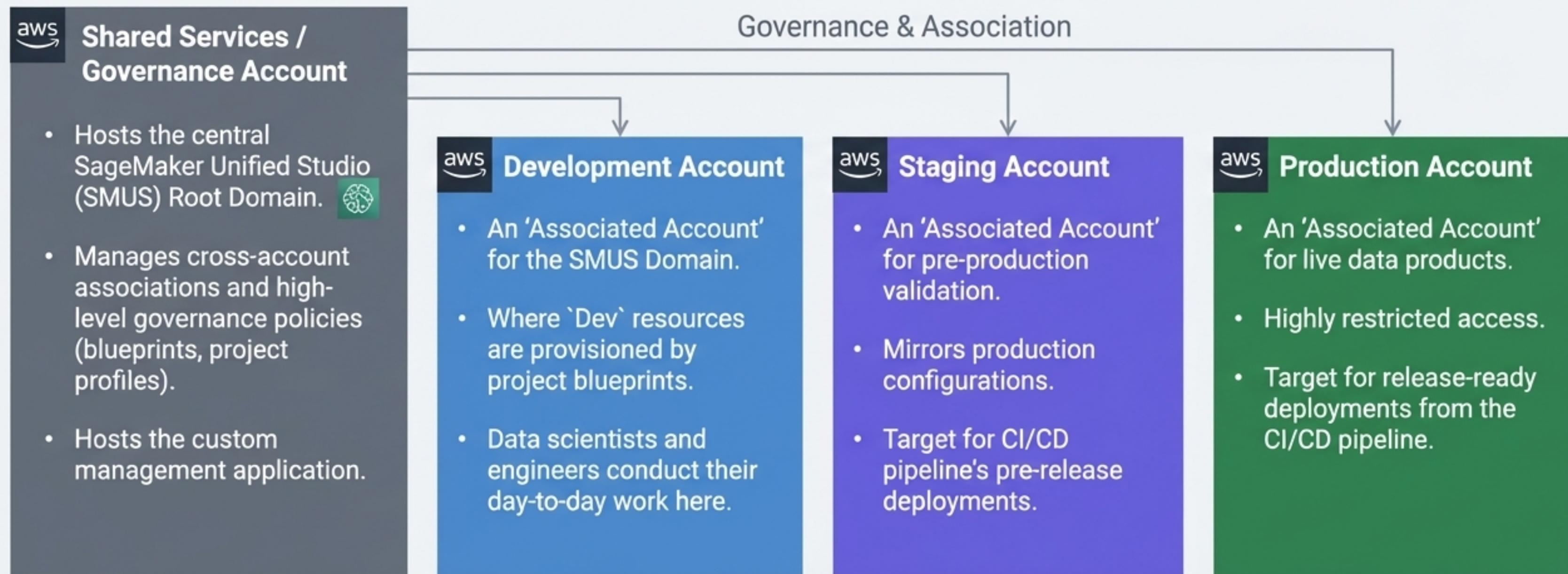
-  **Multi-Account Isolation:** Segregated Dev, Staging, and Production environments in dedicated AWS accounts.
-  **DevOps Automation:** A CI/CD pipeline integrating with an on-premises Git repository.
-  **Programmatic Governance:** A custom application for managing the data mesh (projects, assets, permissions) via AWS SDKs.

Our Approach: This deck follows an **Implementation Journey**. We will build the architecture piece by piece, from the foundational account structure to the operational CI/CD workflow, providing a clear, repeatable pattern.



Foundation: A Four-Account Strategy for Segregation and Control

A robust data mesh begins with a well-defined account structure. We will use a four-account model to enforce strict separation of concerns, enhance security, and simplify governance. The central SMUS Domain will reside in a dedicated Shared Services account.



Translating Data Mesh Concepts to SageMaker Unified Studio

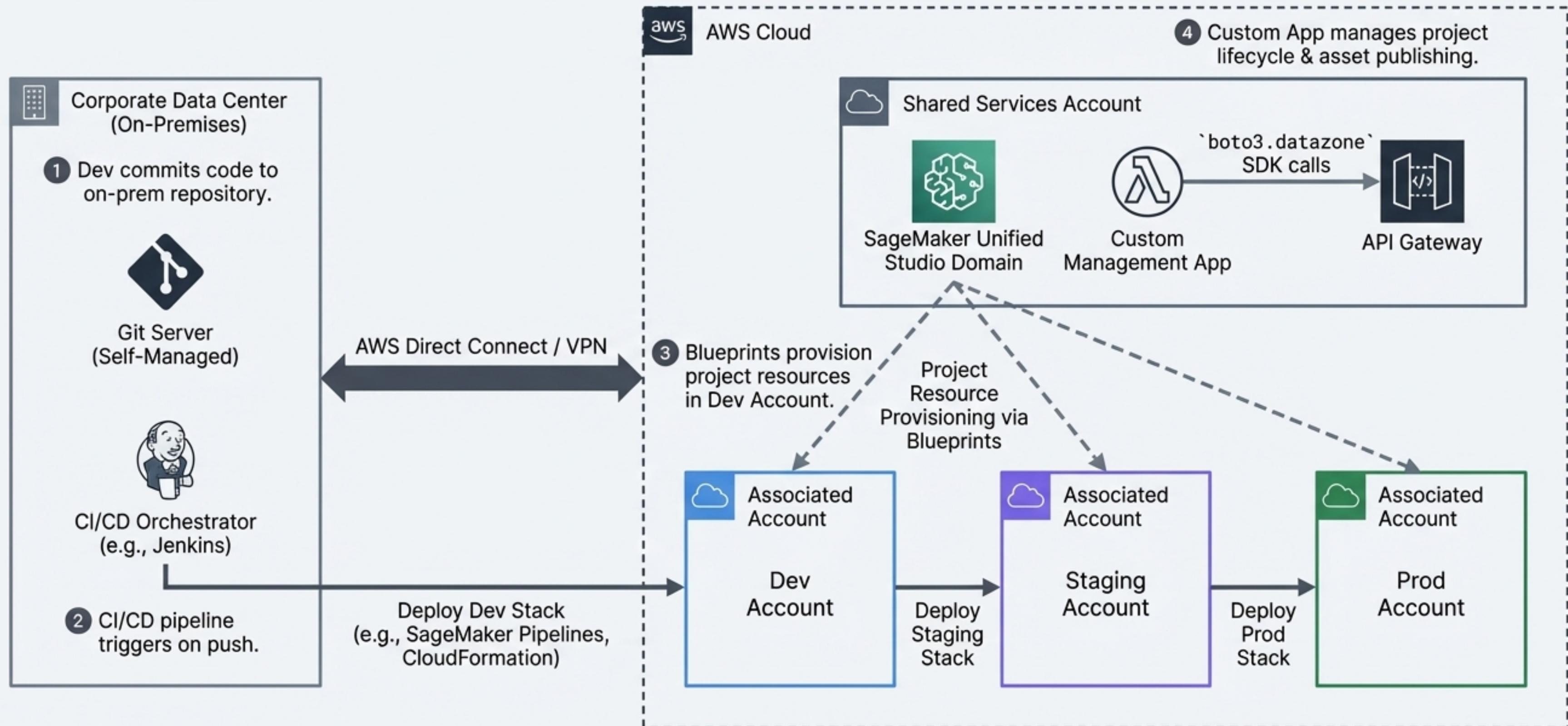
To effectively model our data mesh, we must align our organizational concepts with the core components of SageMaker Unified Studio. This provides a clear and consistent naming convention for our implementation.

Data Mesh Concept	SageMaker Unified Studio (SMUS) Construct	Our Proposed Naming Convention
The Data Mesh Itself	SMUS Domain	A single, top-level domain (e.g., enterprise-data-mesh) in the Shared Services account.
Business Domain (e.g., Marketing, Sales)	SMUS Domain Unit	Hierarchical units that reflect business structure (e.g., MarketingAnalytics, SalesOperations).
Data Product Team / Squad	SMUS Project	The primary workspace for a team, tied to a specific initiative (e.g., CustomerChurnModel-v2).
Central Data Catalog	SMUS Business Data Catalog	The single, domain-wide catalog where all published data assets are discoverable.
Central Business Glossary	SMUS Business Glossary	A collection of business terms associated with assets to ensure consistent definitions across the organization.



Architect's Note: "Think of the **SMUS Domain** as the entire city, **Domain Units** as neighborhoods, and **Projects** as the individual buildings where work gets done. The **Catalog** is the city's public library."

High-Level Architecture: The Complete Picture



Step 1: Establish the Central SMUS Domain

The journey begins by creating the central SMUS Domain in the Shared Services account. This domain acts as the root organizing entity for all users, projects, and assets.

Key Configuration Steps

1. Choose Domain Type:

- **Identity Center-based (Recommended):** Integrates with AWS IAM Identity Center (formerly SSO) for centralized user management. This is the standard for enterprise environments.
- **IAM-based:** Uses traditional IAM roles for authentication. Suitable for simpler setups, but less scalable for large teams.

2. Assign Administrative Permissions:

The IAM role used for setup requires the `SageMakerStudioAdminIAMConsolePolicy` managed policy (rendered in Source Code Pro). This grants the necessary permissions to configure the domain and its resources.

3. Define Domain Execution Role:

An IAM role (`AmazonSageMakerDomainExecutionRole`) must be created to allow SMUS to perform actions on your behalf.

4. Configure Networking & Encryption:

Set up the VPC and subnets where domain resources will reside and choose your data encryption strategy (AWS-managed or customer-managed KMS key).



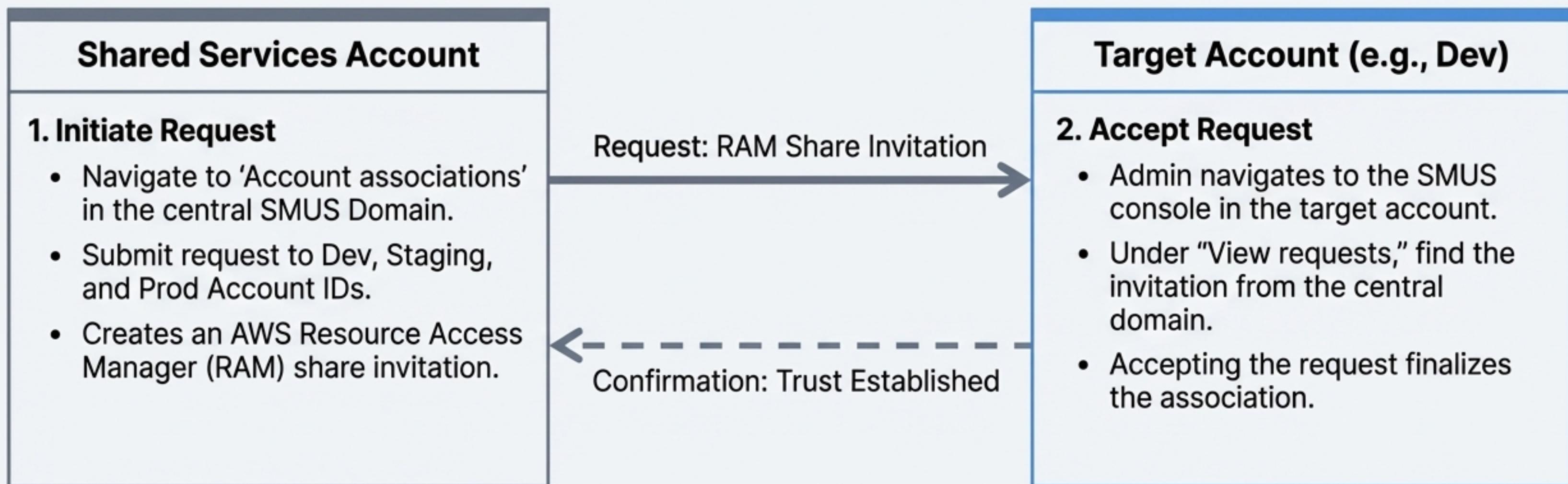
Architect's Note

Per the admin guide, the 'Quick setup' option is intended for testing only. For this production implementation, always use the 'Manual setup' to ensure full control over IAM roles, VPC, and encryption settings.

Step 2: Federate Governance via Associated Accounts

To allow the central SMUS Domain to manage resources across our Dev, Staging, and Prod environments, we must establish a trust relationship using the “Associated Accounts” feature. This is an explicit, bi-directional handshake.

The Association Process

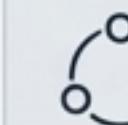
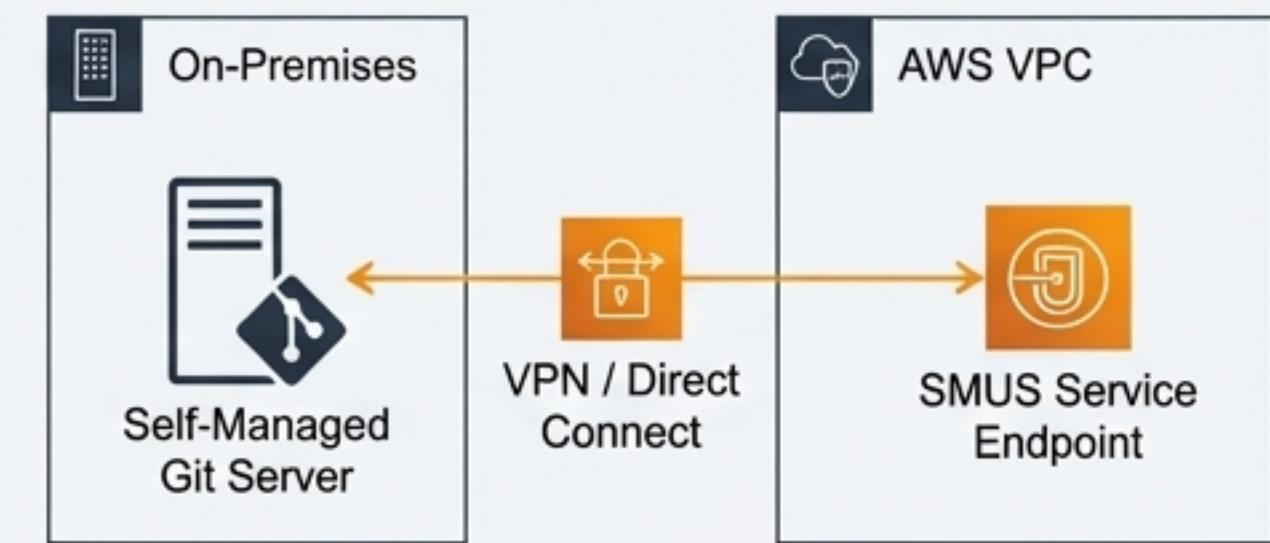


Step 3: Integrate Your On-Premises Git Repository

SMUS provides native support for connecting to self-managed Git servers. This allows our CI/CD pipeline and developer workflows to originate from our existing on-premises source control.

Configuration Process

- 1. Create a New Connection:** In the SMUS Domain settings, under "Connections," select "Create Git connection."
- 2. Select the Correct Provider:** Choose the option that matches your on-prem server: GitHub Enterprise Server or GitLab self-managed.
- 3. Provide Server Details:**
 - URL:** The endpoint of your Git server instance.
 - VPC Configuration:** If the server is not internet-accessible, you must configure the connection to operate within a VPC that has connectivity back to your on-premises network (via Direct Connect or VPN).
 - TLS Certificate:** Provide a custom TLS certificate if required for your server.
- 4. Install the AWS Connector App:** The process involves installing an AWS application on your Git server instance to authorize the connection.
- 5. Enable the Connection:** Once the connection status is "Available," an administrator must explicitly "Enable" it to make it accessible to projects in the domain.



Architect's Note

Network connectivity is paramount. Ensure the VPC security groups and network ACLs allow traffic between the SMUS service endpoints and your on-premises Git server endpoint.

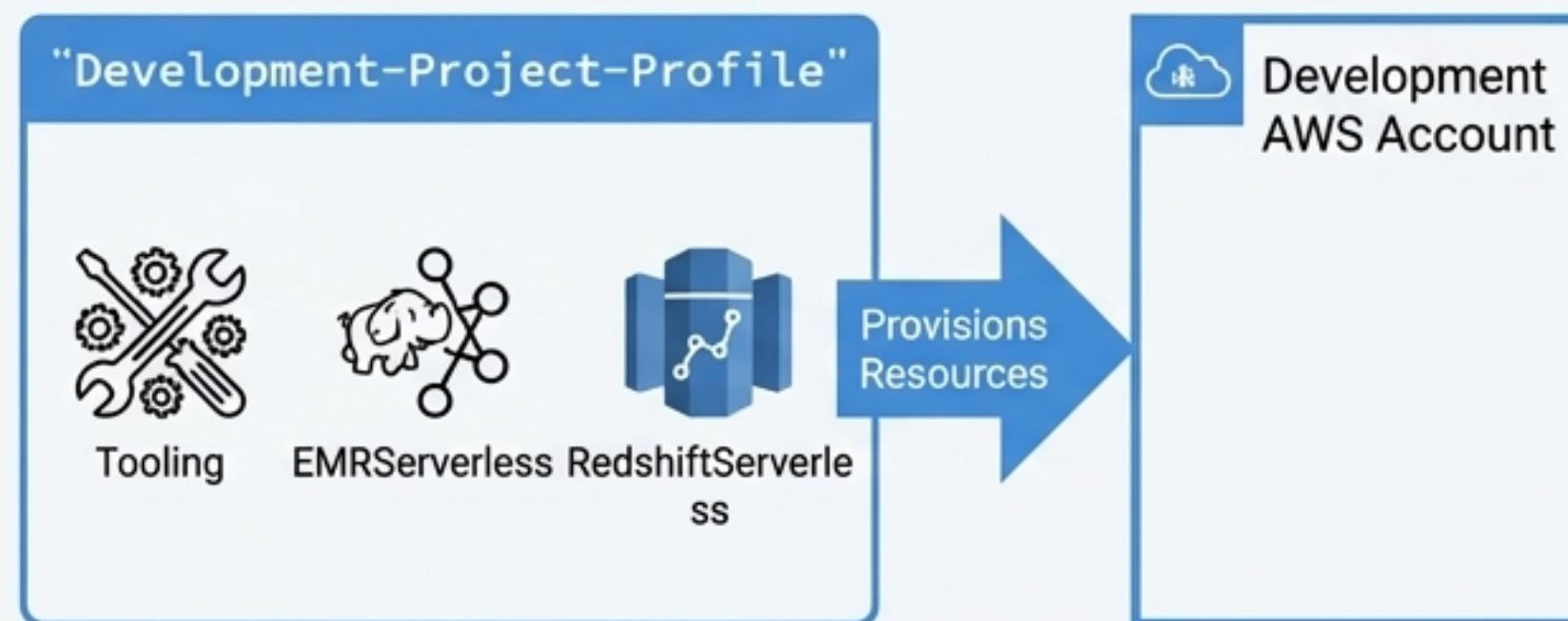
Step 4: Define Environments with Project Profiles and Blueprints

Project Profiles are templates for creating SMUS Projects. They are collections of **Blueprints**, which are CloudFormation templates that provision the necessary AWS resources. We will create distinct profiles for each of our environments.

Our Project Profile Strategy

"Development-Project-Profile"

- Configured to deploy resources into the **Development** associated account.
- Blueprints might include more permissive IAM roles and access to development data sources.
- Can be used by any data scientist to create new projects.



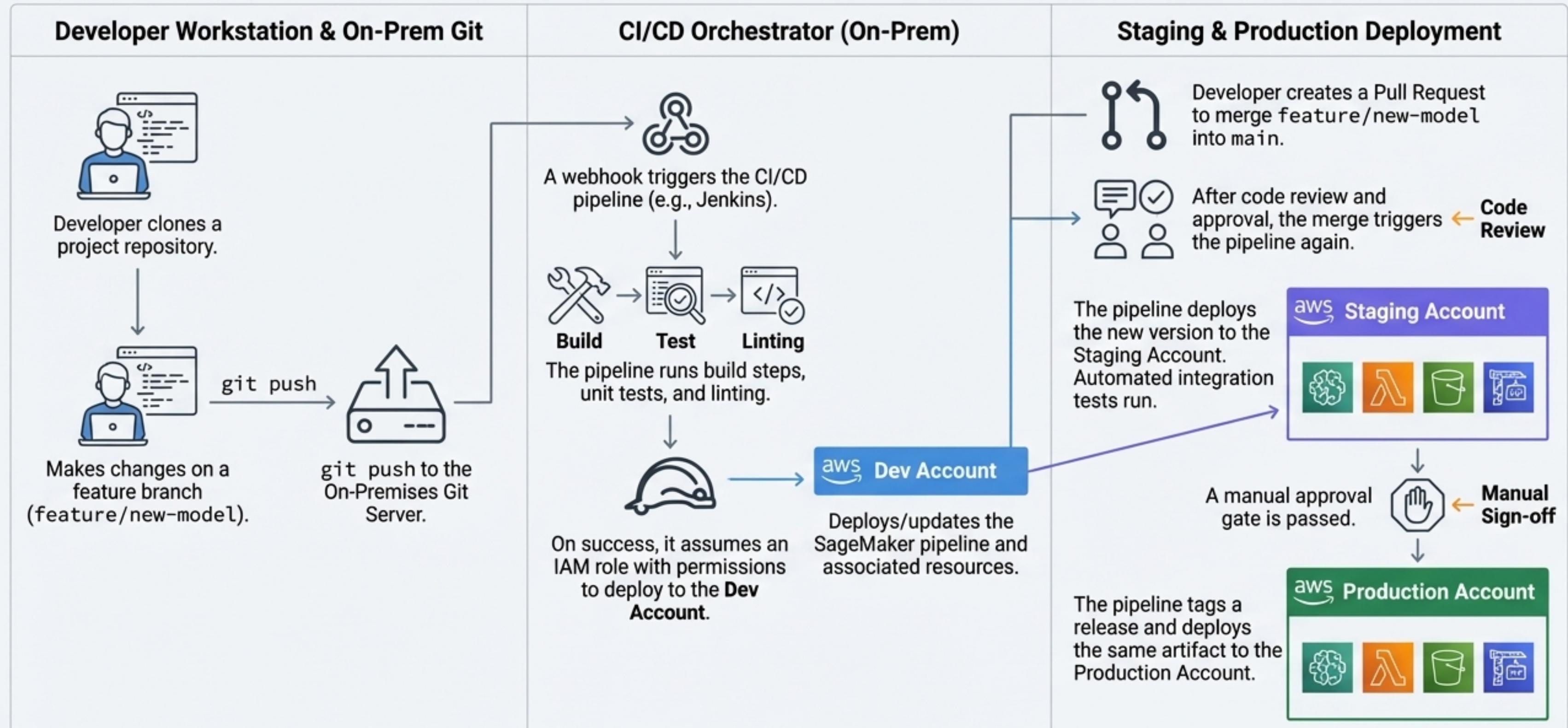
"Staging-Project-Profile"

- Configured to deploy resources into the **Staging** associated account.
- Uses blueprints that mirror production infrastructure and connect to staging data sources.
- Project creation is likely restricted and automated via the CI/CD pipeline.

"Production-Project-Profile"

- Configured to deploy resources into the **Production** associated account.
- Uses hardened, least-privilege blueprints.
- Project creation is strictly controlled and automated.

The CI/CD Workflow in Action



Automating Governance: The Custom Management Application

To manage the data mesh at scale, we will build a custom application that serves as our programmatic control plane. This application will automate routine tasks and enforce governance policies, reducing manual overhead and ensuring consistency. It will be hosted in the Shared Services account.

Key Responsibilities of the Management App

-  **Project Lifecycle Management:** Automate the creation and deletion of SMUs Projects based on team requests.
-  **Asset Management:** Programmatically publish, update, and manage metadata for data assets in the Business Data Catalog.
-  **Subscription Automation:** Handle subscription requests and grants, providing automated access to data products based on predefined rules.
-  **Permissions & Policy Enforcement:** Manage user/group assignments to projects and domain units, ensuring least-privilege access.
-  **Reporting & Auditing:** Generate reports on data asset usage, project activity, and subscription statuses.

Core Technology

The application will be built primarily using the **Boto3 SDK**, specifically the datazone client, which interacts directly with the **SageMaker Unified Studio/Amazon DataZone APIs**.

A Developer's View: Interacting with the SMUS API

The custom management app will orchestrate data mesh operations by calling a suite of AWS APIs via the Boto3 SDK. The `datazone` client is the primary interface for all governance-related tasks.

Key SDK Clients

- `boto3.client('datazone')`: For all core SMUS operations (Domains, Projects, Assets, Subscriptions, Glossaries).
- `boto3.client('sagemaker')`: To manage underlying ML resources like notebooks, training jobs, and endpoints.
- `boto3.client('iam')`: For creating and managing project-specific IAM roles.
- `boto3.client('glue')`,
`boto3.client('lakeformation')`: For interacting with the underlying data sources and permissions.

Pseudo-Code Example: Granting a Data Subscription

```
# Pseudo-code demonstrating a data subscription approval workflow
import boto3

datazone = boto3.client('datazone')
domain_id = 'dzd_enterprise_data_mesh'

def approve_subscription_request(request_id, reason):
    # Retrieve details about the subscription request
    request_details = datazone.get_subscription_request_details(
        domainIdentifier=domain_id,
        identifier=request_id
    )

    # Business logic to validate the request can be added here
    # ...

    # Accept the request, which creates a subscription grant
    response = datazone.accept_subscription_request(
        domainIdentifier=domain_id,
        identifier=request_id,
        acceptMessage={'body': reason}
    )
    print(f"Subscription granted with ID: {response['id']}")
```

Required AWS Services for Implementation

This architecture relies on a set of core AWS services working in concert. Below is a summary of the key components required for the full implementation.

Data Mesh Governance & UI



Amazon SageMaker Unified Studio: The primary interface and governance plane.



Amazon DataZone: The underlying service powering SMUS governance features.



AWS IAM Identity Center: For centralized user authentication and management.

CI/CD & Connectivity



AWS Direct Connect / Site-to-Site VPN: To establish secure connectivity to the on-premises Git server.



Amazon VPC: For network isolation, including subnets, security groups, and endpoints.



(Optional) AWS CodePipeline / CodeBuild: If choosing to host the CI/CD orchestrator on AWS.

Infrastructure & Automation



AWS CloudFormation: The engine used by SMUS Blueprints to provision resources.



AWS IAM: For defining all execution roles, user roles, and service permissions.



AWS KMS: For managing encryption keys for data at rest.

Custom Management Application



AWS Lambda or Amazon ECS/Fargate: To host the application logic.

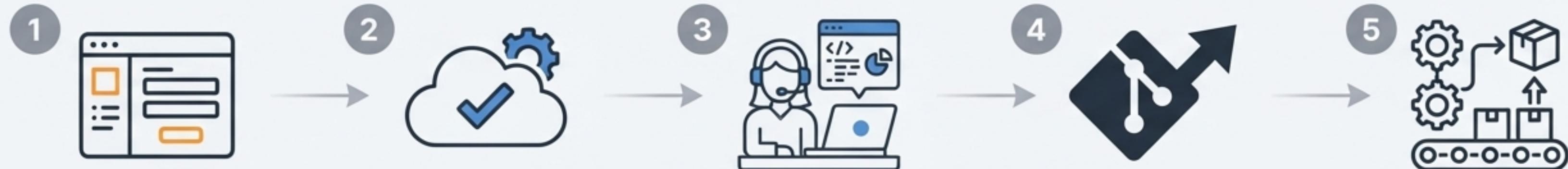


Amazon API Gateway: To expose application functions as a secure API.



Amazon S3: For storing application artifacts and state.

A Day in the Life: The End-to-End Workflow



Project Onboarding

A Team Lead requests a new project via an internal portal. This triggers the Custom Management App. The app calls the `datazone.create_project` API using the Development-Project-Profile.

Environment Provisioning

SMUS automatically creates the project, which includes a dedicated S3 location, IAM roles, and a clone of the appropriate Git repository, all provisioned inside the **Dev Account**. The new Data Scientist is added as a project member.

Development & Experimentation

The Data Scientist accesses the project via the SMUS portal. They launch a JupyterLab space, pull the latest code from the main branch, and create a new feature branch. They write code, query data sources available in Dev, and train a model.

Committing & Pushing

After iterating, they commit their code (`git commit`) and push their feature branch to the On-Prem Git Server (`git push`).

Automated Deployment

The `git push` triggers the CI/CD Pipeline, which automatically deploys their changes to the shared development environment for integration testing, completing the inner development loop.

Your Implementation Roadmap

We recommend a phased approach to build, test, and roll out the data mesh ecosystem. This ensures a solid foundation before scaling to full automation.

Foundational Setup

(Weeks 1-2)

- Provision the four AWS accounts (Shared Services, Dev, Staging, Prod).
- Configure AWS IAM Identity Center and create user groups.
- Deploy the central SageMaker Unified Studio Domain in the Shared Services account.
- Establish network connectivity to the on-premises Git server.

Connectivity & Environment Definition

(Weeks 3-4)

- Configure the Git connection in SMUS for the on-prem server.
- Associate the Dev, Staging, and Prod accounts with the central domain.
- Define and create the initial versions of the [Development](#), [Staging](#), and [Production](#) Project Profiles and their core blueprints.

Automation & Pilot

(Weeks 5-8)

- Build the initial CI/CD pipeline capable of deploying a sample application to the Dev account.
- Develop the Minimum Viable Product (MVP) of the custom management app, focusing on project creation and user assignment.
- Onboard a pilot team to test the end-to-end workflow.