

Git Repo Topology - Data Platform Repos

- Overview
 - Isolation
 - Reusability
 - Complexity
- Topology Options
- What are we actually deploying (and how)?
 - AWS Components
 - Databricks Components
 - Modules
 - Stacks
 - Directories within Stacks
- Rules & Workflow
 - rules

Document Owner(s)	@Athar Shah (Deactivated) @Jay Jiang
Work Package	WP4
Status	REVIEWED
List of Approvers	<div><div><input checked="" type="checkbox"/> @Sampath Jagannathan / @Vishnu Devarajan (Unlicensed)</div><div><input checked="" type="checkbox"/> @Neil Belford</div><div><input checked="" type="checkbox"/> @Former user (Deleted) / @Brandon Lay (Unlicensed) / @Davood Shaiek</div><div><input checked="" type="checkbox"/> @Wee Lih Lee / @Nikilesh Chivukula / @Lawrence Law</div></div>

Overview

This data platform repo design seeks to strike a balance between 3 competing principles.

Isolation

Environments and BUs should be isolated from each other to prevent collisions. We want to ensure that a change in a development environment does not inadvertently force a refresh of the production state (same applies for changes in other BUs).

Reusability

Also known as DRY (Don't Repeat Yourself); platform code is often used to deploy for multiple instances of the same resource with minor variations. We will use [terraform modules](#) to abstract the common patterns across BUs to maximise code re-usability. This principle often conflicts with **isolation** principle, as making code re-usable also create dependencies between **modules** and callers of **modules** ; i.e. if 2 BUs share a common module, and BU-1 requires additional functionality added to a shared **module** , changes to the shared **module** must **either** be regression tested to ensure BU-2 will not be impacted **or** the **module** needs to have semantic version tags to adhere to **isolation** principle.

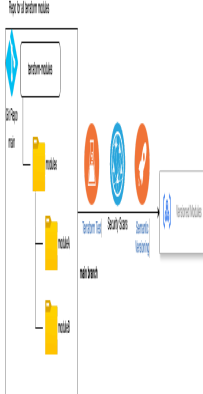
Complexity

The conflict between **isolation** and **reusability** often result in added **complexity** . In the previous example, having semantic versions enabled on the **modules** often mean that you need to break the modules out into their own repos. Considering the size of the operational team, we want to minimise the number of repos for Jemena (as typically more repos means more complexity); however, this means that one repo will be able to have multiple terraform statefiles (corresponding to subdirectories in the repo). The one-to-many relationship between repo and states can be seen as complexity, however, this is the type of complexity that the majority of CI/CD patterns are designed to handle.

Topology Options

Infrastructure Git Topology Options

Separate Terraform Modules Report



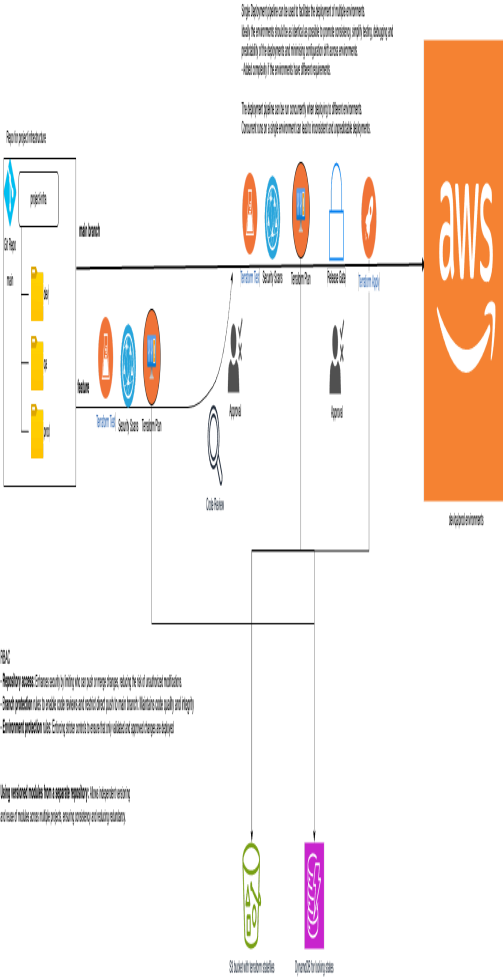
Allows independent versioning, and reuse of modules across multiple projects, ensuring consistency and reducing redundancy

Modules can be pinned in the Infrastructure code by the `branch` command.

Complexity	
	1 Government assembly: single-appeal department 2 Government assembly: multiple-appeal departments 3 Government assembly: province-appeal department 4 Government assembly: province-multiple departments
ESG	ESG operational in multiple departments ESG operational in multiple departments ESG operational in multiple departments ESG operational in multiple departments
Market diversity	Supports various principles and standards for market diversity Supports various principles and standards for market diversity Supports various principles and standards for market diversity Supports various principles and standards for market diversity
Policy uptake	Get environment-friendly policies enacted as part of the public Get environment-friendly policies enacted as part of the public Get environment-friendly policies enacted as part of the public Get environment-friendly policies enacted as part of the public
Search capacity	Add up government's search capacity to environment Add up government's search capacity to environment Add up government's search capacity to environment Add up government's search capacity to environment

Approach 1: Separate Folder for each environment - Single Deployment Pipeline

Trunk based approach



to use this state before using `DynamicIO` to prevent concurrent operations

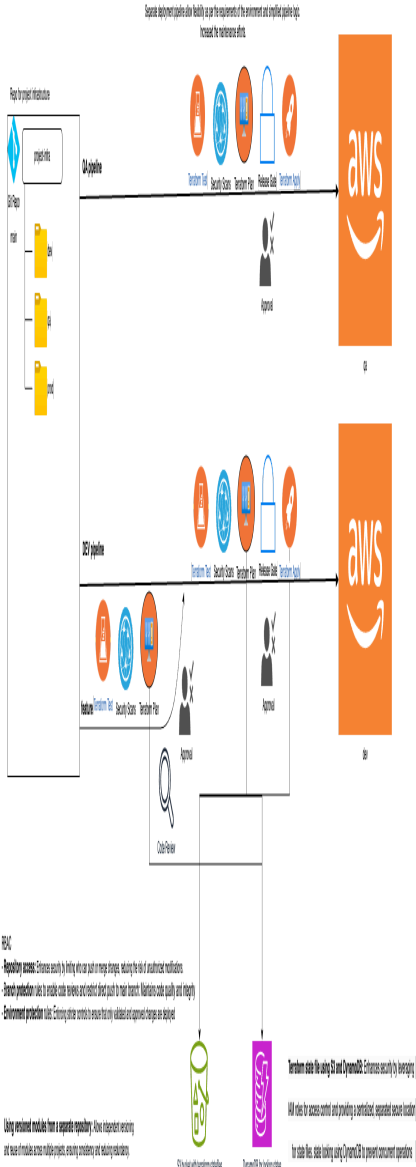
Approach 3: Long Running Branches for each environment - Single Deployment Pipeline

Gilow: Long running branch



Approach 2: Separate Folder for each environment - Separate Deployment Pipeline

Trunk based accounts



Approach 4: Long Running Branches for each environment - Separate Deployment Pipeline

Gifts: Long running branch



core-network- databricks- vpc- components	Private Link ENI, Security Groups	Jemena / Projects / Future Networks Datahub / Databricks / core-network-databricks-vpc-components · GitLab
app-datahub- nonprod- databricks- aws-infra	S3, Instance Profiles, IAM Roles	Jemena / Projects / Future Networks Datahub / Databricks / app-datahub-nonprod-databricks-aws-infra · GitLab
app-datahub- prod- databricks- aws-infra	S3, Instance Profiles, IAM Roles	Jemena / Projects / Future Networks Datahub / Databricks / app-datahub-prod-databricks-aws-infra · GitLab

Databricks Components

After AWS components are deployed, they either need to be registered in databricks or otherwise used by databricks to create subsequent resources, in order to pass AWS ARNs etc. onto databricks resources, we will need to use GitLab's [terraform remote state as a data source](#) and provision the corresponding [permission for CI runners](#).

Repo	Resource Level	Link
databricks- unity-catalog	Account Only + Workspace-agnostic resources	Jemena / Projects / Future Networks Datahub / Databricks / databricks-unity-catalog · GitLab

databricks- workspaces	Account + Workspace	Jemena / Projects / Future Networks Datahub / Databricks / databricks-workspaces · GitLab
---------------------------	---------------------	---

Within these 2 repos, the first layer of directories distinguish **modules** and **stacks**.

Modules

Modules sub-directory contain the definition for how to deploy a particular combination of resources. e.g. **RBAC** module would define how the RBAC hierarchy is defined for any arbitrary Business Unit. Nothing is created by modules.

Stacks

Stacks call **modules**, they pass in configuration to modules and create the resources based on the patterns defined within modules. e.g. **RBAC** stack would call **RBAC** module (via source=) and pass in configuration such as BU's name.

Directories within Stacks

Every directory within stacks directory has its own terraform statefile, this is how we ensure **isolation** across BUs. Environment injection happens dynamically via

x.auto.tfvars within those sub directories to distinguish **prod** and **noprod**

deployments; the injection of actual values into **x.auto.tfvars** will be set based on the CI/CD design (typically split by branch names).

```

1  |─ databricks-unity-catalog
2    |   |─ modules
3    |   |   |─ rbac
4    |   |   |─ schemas
5    |   |   |─ service-principals
6    |   |─ stacks
7    |   |   |─ digital-bu
8    |   |   |   |─ catalogs.tf
9    |   |   |   |─ rbac.tf
10   |   |   |   |─ schemas
11   |   |   |─ elec-network-bu
12   |   |   |   |─ catalogs.tf
13   |   |   |   |─ rbac.tf
14   |   |   |   |─ schemas
15   |   |   |─ org (no CI, manual deploy)
16   |   |   |   |─ metastore.tf
17   |─ databricks-workspaces

```

```
18 | modules
19 |   ├── cluster-policy
20 |   ├── long-lived-clusters
21 |   ├── managed-workspace
22 |   └── sql-warehouses
23 | stacks
24 |   ├── digital-lab-workspace-infra
25 |   ├── digital-field-workspace-infra
26 |   ├── digital-lab-workspace-resources
27 |   ├── digital-field-workspace-resources
28 |   ├── elec-network-lab-workspace-infra
29 |   ├── elec-network-field-workspace-infra
30 |   ├── elec-network-lab-workspace-resources
31 |   └── elec-network-field-workspace-resources
```

Rules & Workflow

rules

- main direct push to **main** denied, can only be modified via **MR**
- **MR**s require minimum 1 approval and pipeline success to merge
 - subsequent commits revoke previous approvals