

1.2 AIB2.0 User Journey Workflows

- [1. Background & Context](#)
- [2. Key Git Repositories](#)
- [3. Key User Journeys](#)

1. Background & Context

AIBooster 2.0 aims to build a generic GenOps Platform on top of the existing AIB1.0 MLOps Platform while providing functional systems & templates to rapidly productionize a Generative AI business use case while maintaining best practices and encapsulating all possible aspects that an enterprise-grade platform should capture (scalability, robustness, logging, monitoring etc)

A developer leveraging this platform can adhere to multiple personas like an AIB ML Engineer, Use Case ML Engineer, AIB Platform Engineer etc.

This playbook aims to guide developers new to the platform through key user journeys to get started with the same in a simplified & concise manner.

#	Git Repo	Description	Expected end Outcome
		<p>requirements while ensuring minimal impact to the process of productionizing a GenAI use case</p> <p>This service is deployed as a Cloud Run instance in an AI Gateway Orchestration layer</p> <p>For more details, please refer to 1.3.2 Responsible AI Service</p>	the same can be used as a base to deploy the service to Cloud Run
2	LLM Proxy Service	<p>This repository covers the logic running for the LLM Proxy Service.</p> <p>Any LLM interactions that a GenAI use case needs to run passes through this service providing a framework to control which LLMs can a use case access & what kind of Responsible AI configuration is required for the same</p> <p>For more details, please refer to 1.4.1 LLM Proxy Service</p>	<p>This repository outputs an image stored in a centrally located Google Cloud Artifact Registry via CI/CD &</p> <p>the same can be used as a base to deploy the service to Cloud Run</p>
3	AIB GenAI Template	<p>This repository provides a reference implementation for productionizing a GenAI use case on Google Cloud, & can be used as a starting point w/ configurable boilerplate code. The implementation includes:</p> <ul style="list-style-type: none"> • Data ingestion and vectorisation pipeline using Kubeflow Pipelines and Vertex AI Pipelines. • LLM inference using Vertex AI with RAG setup by Vertex AI Vector Search. • Infrastructure-as-Code using Terraform for Cloud Run services, PubSub and Cloud Scheduler. 	<p>This repository once cloned for a specific use case, can output:</p> <ul style="list-style-type: none"> • A Base image to be used by the overall use case logic (RAG Inference + Data Vectorisation + LLM Evaluation) • A Cloud Scheduler to run Data Vectorisation regularly • A Cloud Scheduler to run Evaluation pipelines regularly
4	AIB Inference Deployments	<p>This repository serves as a centrally-managed place to deploy any GenAI Use Case application as a Cloud Run Service.</p> <p>Using YAML-based configurations hooked w/ automated CI/CD under-the-hood, a developer can quickly plug in configuration parameters to deploy an application across environments (lab/non-live/live)</p>	This repository runs terraform-based deployment of all GenAI Use Case applications to Cloud Run via automated CI/CD under the hood

3. Key User Journeys

Key User Journeys summarized below capture how different personas can leverage the 4 Git Repos mentioned above to successfully complete a relevant user journey

#		Persona		User Journey	Applicable Git Repo
1	As a/an	AIB ML Engineer	,I want to	Update the Responsible AI git repo	Responsible AI Service
2				Add a new Responsible AI module	
3				Deploy a new Responsible AI service image version	
4				Update the Proxy service git repo	LLM Proxy Service
5				Add a new Proxy router	
6				Deploy a new Proxy service image version	
7		Use Case ML Engineer/ Data Scientist		Build an image for a GenAI use case	AIB GenAI Template
8				Develop & test a Data Vectorisation pipeline	
9				Create a Data Vectorisation pipeline schedule & Promote to higher environments	
10				Develop & test the RAG inference logic	AIB Inference Deployments
11				Deploy a RAG inference use case application & Promote to higher environments	
12				Develop & Test an Evaluation pipeline - Golden Dataset pattern	AIB GenAI Template
13				Develop & Test an Evaluation pipeline - Continuous evaluation pattern	
14				Create an Evaluation pipeline schedule & Promote to higher environments	
15				Update Proxy config for my use case	

Workflows for these 15 User Journeys can be summarized in the designs below w/ further details on each User Journey captured in subsequent sub-pages:

# Element	Description	Assumptions
		using this guide - How to : Connect to AI Booster Notebook vis SSH - AI Booster
2 Add a python file for the new scanner	<p>All Responsible AI modules exist in the <code>/responsible_ai_modules/scanners/</code> folder.</p> <p>This includes existing modules for:</p> <ul style="list-style-type: none"> • Content Moderation using the Google Cloud Natural Language API • Anonymization using the Google Cloud DLP API • Prompt Injection using an open-source model deployed to a Vertex endpoint in the AI Gateway layer <p>New modules can include future roadmap functionality like topicality, embedding scans etc.</p> <p>Any new module needs to exist in the above-mentioned folder w/ a base python class structure containing:</p> <ul style="list-style-type: none"> • The new scanner class deriving from the Base Scanner class • 2 methods for: <ul style="list-style-type: none"> ○ base initialisation of the class ○ Scan method capturing underlying logic for the new module <p>Furthermore, all scanners/modules are generalized to adhere to a pydantic model called <code>IndividualScannerOutput</code> that exists in <code>responsible_ai_modules/schemas.py</code>. Any new module needs to transform their output to adhere to this structure</p>	
3 Update supplementary files to include this new scanner	<p>In addition to a new scanner file, following supplementary files need to be updated to include this new module:</p> <ul style="list-style-type: none"> • <code>responsible_ai_modules/scanners/__init__.py</code> • <code>get_scanner_by_name</code> function in <code>responsible_ai_modules/scanners/utils.py</code> 	
4 Add unit tests	The repo contains unit tests related to all existing functionality in the <code>tests/</code> folder	

# Element	Description	Assumptions
	This markdown file includes further details on testing the service locally & miscellaneous related information on unit testing etc	Notebook - AI Booster <ul style="list-style-type: none"> VSCode has been connected to the created Workbench instance using this guide - How to : Connect to AI Booster Notebook vis SSH - AI Booster
2 Refactor code as required	Once the development workspace has been completely setup, a developer can create a feature branch, refactor code as required & commit changes to said feature branch	
3 Update existing unit tests	<p>The repo contains unit tests related to all existing functionality in the <code>tests/</code> folder</p> <p>These tests can be updated based on the refactoring of functionality as required & tested locally by executing the <code>make pytest</code> command.</p> <p>This Makefile abstracts away execution of key commands & can be referred for other aspects as well</p>	
4 Test changes locally	<p>Once required changes have been committed, the same can be tested locally by spinning up the LLM Proxy Service locally using <code>make run</code> which spins up the FastAPI server on localhost:8000 by default</p> <p>Local testing can be executed in 2 ways:</p> <ul style="list-style-type: none"> <code>curl</code> directly into the local server Spin up a dummy use case app locally using the <code>make run-dummy-use-case</code> command & update the <code>dummy_use_case/api/app.py</code> file as required <p>The 2nd option provides more flexibility to test any LLM Proxy changes from an overall usage standpoint (calling the LLM Proxy via Langchain, Llama-Index or other LLM frameworks)</p>	
5 Create a Pull request	Once relevant changes have been tested & the feature branch is good to go, next logical step would be to create a Pull Request in the repo for review	
6 Merge to main	Post review & approval of the feature branch functionality, the same can be merged into the <code>main</code> branch	

# Element	Description	Assumptions
	<p>New routers can include LLMs that do not conform to the extensibility provided by above-mentioned routers. In that scenario, a new router file needs to be added to this folder w/ relevant authentication logic & API calling logic.</p> <p>Please refer to existing Google LLM, External LLM routers for ideas on standardized code structure</p> <p>Once logic for the new router has been plugged in, this router needs to be included in the base <code>api/app.py</code> file</p>	
3 Add logic for parsing relevant new model patterns	<p>Every LLM call via a router conforms to a generalized regex model pattern that captures all possible model paths as per the final LLM endpoint. This is currently captured in the <code>env.sh</code> file & the same is deployed with LLM Proxy Cloud Run deployment as environment variables.</p> <p>In addition to the above, 2 utility files capture logic to generically parse these model patterns & the same need to be amended:</p> <ul style="list-style-type: none"> • api/utls/model_path_parsing.py - Need to add a new case for the new model patterns • api/utls/schema.py - Need to update the <code>ModelProvider</code> class to reflect this new model pattern 	
4 Add unit tests	<p>The repo contains unit tests related to all existing functionality in the <code>tests/</code> folder</p> <p>This folder will need to be updated w/ new unit tests to reflect new functionality created for the new scanner.</p> <p>This can be tested locally by executing the <code>make pytest</code> command.</p> <p>This Makefile abstracts away execution of key commands & can be referred for other aspects as well</p>	
5 Test changes locally	<p>Once required changes have been committed, the same can be tested locally by spinning up the LLM Proxy Service locally using <code>make run</code> which spins up the FastAPI server on localhost:8000 by default</p> <p>Local testing can be executed in 2 ways:</p> <ul style="list-style-type: none"> • <code>curl</code> directly into the local server • Spin up a dummy use case app locally using the <code>make run-dummy-use-case</code> command & update the <code>dummy_use_case/api/app.py</code> file as required 	

#	Element	Description	Assumptions
3a	CI process automatically builds & pushes image to Dev Booster Mirror	<p>Once a tag w/ dev as suffix is created, the CI pipeline is automatically triggered that runs the image build & push process & finally outputs an image in Google Cloud Artifact registry of Dev Booster Mirror</p> <p>For more details on this CI pipeline run by Cloud Build, please refer to the <code>Dockerfile</code> & <code>cloudbuild.yaml</code> files in the repo</p>	<p>Cloud Build Trigger exists in the Dev Booster Mirror Google Cloud project to run this CI pipeline.</p> <p>This is a manual one-time activity & is already in place</p>
3b	CI process automatically builds & pushes image to Prod Booster Mirror	<p>Once a tag w/ prd as suffix is created, the CI pipeline is automatically triggered that runs the image build & push process & finally outputs an image in Google Cloud Artifact registry of Prod Booster Mirror</p> <p>For more details on this CI pipeline run by Cloud Build, please refer to the <code>Dockerfile</code> & <code>cloudbuild.yaml</code> files in the repo</p>	<p>Cloud Build Trigger exists in the Prod Booster Mirror Google Cloud project to run this CI pipeline.</p> <p>This is a manual one-time activity & is already in place</p>

1.2.2 User Journeys - Persona: Use Case ML Engineer/ Data Scientist

- [1. Persona Description](#)
- [2. Persona-specific User Journeys](#)
- [2.1. User Journey #7 - Build an image for a GenAI use case](#)
- [2.2. User Journey #8 - Develop & Test a Data Vectorisation Pipeline](#)
- [2.3. User Journey #9 - Create a Data Vectorisation Pipeline Schedule & Promote to higher environments](#)
- [2.4. User Journey #10 - Develop & Test the RAG Inference logic](#)
- [2.5. User Journey #11 - Deploy a RAG Inference Use Case Application & Promote to higher environments](#)
- [2.6. User Journey #12 - Develop & Test an Evaluation Pipeline - Golden Dataset pattern](#)
- [2.7. User Journey #13 - Develop & Test an Evaluation Pipeline - Continuous Evaluation pattern](#)
- [2.8. User Journey #14 - Create an Evaluation Pipeline Schedule & Promote to higher environments](#)

1. Persona Description

As a Use Case ML Engineer or a Use Case Data Scientist, key responsibilities can include leveraging the AIB GenAI Template as boilerplate code to build out relevant logic for a specific GenAI use case & productionize the same w/ components provided by the AIB2.0 Platform.

These templates capture the following essential elements of a GenAI use case:

#	Element	Description	Assumptions
		<p>as required & commit changes to said feature branch</p> <p>The repo also contains unit tests related to all existing functionality in the <code>testing/</code> folder</p> <p>These tests can be updated based on the refactoring as required & tested locally by executing the <code>make pytest</code> command.</p> <p>This Makefile abstracts away execution of key commands & can be referred for other aspects as well</p>	
4	Create a pull request	Once relevant changes have been tested & the feature branch is good to go, next logical step would be to create a Pull Request in the repo for review	
5a	Merge to main	Post review & approval of the feature branch functionality, the same can be merged into the <code>main</code> branch	
5b	Create a Release tag	<p>A new release tag needs to be created w/ relevant description supporting it.</p> <p>This tag should conform to best practices around tags, specifically of format:</p> <p>vMAJOR.MINOR.PATCH</p> <p>eg - A new release tag can be:</p> <ul style="list-style-type: none"> v1.2.0 	
6a	CI process automatically builds & pushes an image to the build-nl Use Case Vendor layer project	<p>Once the PR has been pushed to the <code>main</code> branch (see step 5a), the CI/CD pipeline is automatically triggered that:</p> <ul style="list-style-type: none"> Runs relevant unit testing & linting checks Runs the image build & push process & finally outputs an image in Google Cloud Artifact registry of the Use Case Vendor layer build-nl project with a <code>latest</code> tag along with the relevant git commit SHA value Compiles relevant Vertex pipelines & stores them in GCS (more details to follow in subsequent user journeys) Create/update a Cloud Scheduler job for the compiled Vertex pipeline (more details to follow in subsequent user journeys) 	<p>Cloud Build Trigger exists in the Use Case Vendor layer build-nl project to run this CI/CD pipeline</p> <p>This is a one-time manual activity & needs to be implemented once for every GenAI use case</p>

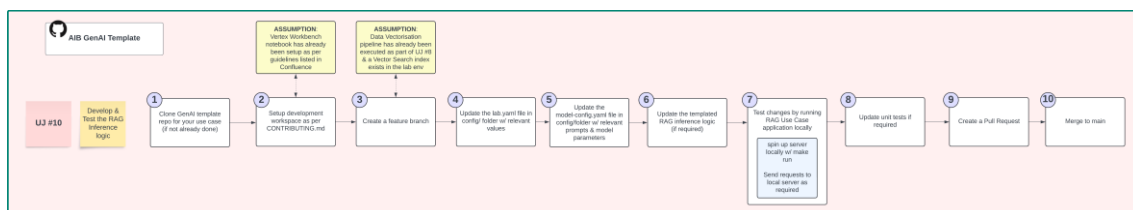
#	Element	Description	Assumptions
			Notebook - AI Booster <ul style="list-style-type: none"> VSCode has been connected to the created Workbench instance using this guide - How to : Connect to AI Booster Notebook vis SSH - AI Booster
3	Create a feature branch	Create a feature branch on the cloned repo to build & test relevant Data Vectorisation pipeline functionality	
4	Update the lab.yaml file in config/ w/ relevant values	<p>In addition to boilerplate code to build a Vertex pipeline for Data Vectorisation & ingest documents to a Vector Database (Vertex Vector Search), this template abstracts relevant updates to a configuration-based YAML file which exists in the config/ folder</p> <p>This config folder will need a sub-folder w/ naming of the AIB instance (eg vf-uk-aib-prd-cip-ask-lab/nl/live points to vf-uk-aib-prd-cip-ask) & YAML files within said sub-folder pointing to the lab, non-live & live instances</p> <p>For testing purposes, the lab.yaml file needs to be updated w/ relevant values catering to:</p> <ul style="list-style-type: none"> Project ID Region Relevant Google Cloud Storage URIs Relevant Vertex Pipeline configuration Vertex Vector Search fields <p>The template repo contains sample YAML files for reference</p>	Source data files to be ingested into a Vector Database already exist in a relevant Google Cloud Storage bucket in the relevant Use Case Vendor layer project
5	Update code for Vertex pipeline components & the pipeline (if required)	<p>Once the configuration-based YAML files have been updated, base logic in the <code>ingestion_workflow/pipeline/</code> folder can be updated if absolutely required.</p> <p>This folder contains generic logic for Vertex Pipeline components that include:</p> <ul style="list-style-type: none"> Fetching documents Processing documents 	

#	Element	Description	Assumptions
		<ul style="list-style-type: none"> Creating or Updating a Vertex Vector Search index (Vector Database) Updating pipeline details <p>All components above contain core logic to maintain & audit the document lifecycle of processed documents</p>	
6	Update unit tests (if required)	<p>The repo contains unit tests related to all existing functionality in the <code>testing/</code> folder</p> <p>These tests can be updated based on the refactoring as required & tested locally by executing the <code>make pytest</code> command.</p> <p>This Makefile abstracts away execution of key commands & can be referred for other aspects as well</p>	
7	Manually compile & submit a Vertex Pipeline job	<p>Once all changes are in place, the Vertex pipeline can be manually submitted via a single python command -</p> <pre>python ingestion_workflow/pipeline/pipeline_run.py</pre> <p>to test updated functionality</p>	
8a	Validate Vertex Pipeline execution	<p>Once the above-mentioned python command has been executed in a development workspace, the Terminal will provide a link for the running pipeline OR the same can be found in the Vertex Pipeline UI</p> <p>Once submitted & executed, this pipeline should be validated via the UI to verify updated functionality & configuration (if any)</p>	
8b	Validate create Vector Database	<p>Once the Vertex pipeline has been submitted & successfully executed, a Vertex Vector Search index will be created/updated w/ embeddings of the ingested documents.</p> <p>The same should be validated from a use case specific standpoint</p>	
9	Create a Pull Request	<p>Once relevant changes have been tested & the feature branch is good to go, next logical step would be to create a Pull Request in the repo for review</p>	
10	Merge to main	<p>Post review & approval of the feature branch functionality, the same can be merged into the <code>main</code> branch</p>	

#	Element	Description	Assumptions
3	Update the cron schedule for the Cloud Scheduler job	In addition to relevant configuration-based YAML files, separate files contain logic to configure the cron schedule of the Cloud Scheduler job to be created/updated cron schedules for both non-live & live Use Case Vendor layer projects will need to be updated separately	
4	Create a Pull Request	Once relevant changes have been tested & the feature branch is good to go, next logical step would be to create a Pull Request in the repo for review	
5a	Merge to main	Post review & approval of the config-update/feature branch functionality, the same can be merged into the <code>main</code> branch	
5b	Create a Release tag	A new release tag needs to be created w/ relevant description supporting it. This tag should conform to best practices around tags, specifically of format: <code>vMAJOR.MINOR.PATCH</code> eg - A new release tag can be: <ul style="list-style-type: none"> • <code>v1.2.0</code> 	
6a	CI process automatically creates/updates the Cloud Scheduler job in the non-live Use Case Vendor layer project	Once the PR has been pushed to the <code>main</code> branch (see step 5a), the CI/CD pipeline is automatically triggered that: <ul style="list-style-type: none"> • Runs relevant unit testing & linting checks • Runs the image build & push process & finally outputs an image in Google Cloud Artifact registry of the Use Case Vendor layer build-nl project with a <code>latest</code> tag along with the relevant git commit SHA value • Compiles relevant Vertex pipelines & stores them in GCS under a folder w/ name mapping to the git commit SHA value in the build-nl Use Case Vendor layer project • Create/update a Cloud Scheduler job for the compiled Vertex pipeline in the non-live Use Case Vendor layer project <ul style="list-style-type: none"> ○ This is done via Terraform embedded into the CI/CD pipeline as its final stage & it utilizes the 	Cloud Build Trigger exists in the Use Case Vendor layer build-nl project to run this CI/CD pipeline This is a one-time manual activity & needs to be implemented once for every GenAI use case

#	Element	Description	Assumptions
		git commit SHA value attached to the image + the compiled pipeline GCS location from the previous stage	
6b	CI process automatically creates/updates the Cloud Scheduler job in the live Use Case Vendor layer project	<p>Once a new release tag has been created (see step 5b), the CI/CD pipeline is automatically triggered that:</p> <ul style="list-style-type: none"> Runs relevant unit testing & linting checks Runs the image build & push process & finally outputs an image in Google Cloud Artifact registry of the Use Case Vendor layer build-live project with the release tag attached to the image Compiles relevant Vertex pipelines & stores them in GCS under a folder w/ name mapping to the release tag in the build-live Use Case Vendor layer project Create/update a Cloud Scheduler job for the compiled Vertex pipeline in the live Use Case Vendor layer project <ul style="list-style-type: none"> This is done via Terraform embedded into the CI/CD pipeline as its final stage & it utilizes the release tag attached to the image + the compiled pipeline GCS location from the previous stage 	<p>Cloud Build Trigger exists in the Use Case Vendor layer build-live project to run this CI/CD pipeline</p> <p>This is a one-time manual activity & needs to be implemented once for every GenAI use case</p>

2.4. User Journey #10 - Develop & Test the RAG Inference logic



Based on the user journey outlined above, this UJ applies to the [AIB GenAI Template git repo](#) & details for the same can be broken down as follows:

#	Element	Description	Assumptions
1	Clone GenAI Template repo	Clone the template repo for your GenAI use case & leverage the boilerplate code & frameworks already provided by the template	

#	Element	Description	Assumptions
2	Setup a development workspace	<p>Once the assumptions have been met, further setup is required. This includes:</p> <ul style="list-style-type: none"> • Python package installation using poetry as a package manager • Loading environment variables <p>Detailed steps for the same are captured in the CONTRIBUTING.md of the repo</p>	<p>Development workspace has been setup by:</p> <ul style="list-style-type: none"> • Creating a Vertex Workbench instance using this guide - How to : Create a Vertex Notebook - AI Booster • VSCode has been connected to the created Workbench instance using this guide - How to : Connect to AI Booster Notebook vis SSH - AI Booster
3	Create a feature branch	Create a feature branch on the cloned repo to build & test relevant RAG inference logic	<p>Data Vectorisation pipeline has already been successfully executed as part of UJ #8 &</p> <p>a Vector Search index exists containing embeddings of the ingested documents</p>
4	Update the lab.yaml file in config/ w/ relevant values	<p>In addition to boilerplate FastAPI code to build out RAG inference logic, this template abstracts relevant updates to a configuration-based YAML file which exists in the config/ folder</p> <p>This config folder will need a sub-folder w/ naming of the AIB instance (eg vf-uk-aib-prd-cip-ask-lab/nl/live points to vf-uk-aib-prd-cip-ask) & YAML files within said sub-folder pointing to the lab, non-live & live instances</p> <p>For testing purposes, the lab.yaml file needs to be updated w/ relevant values catering to:</p> <ul style="list-style-type: none"> • Project ID • Region • Relevant Google Cloud Storage URIs • Relevant Vertex Pipeline configuration • Vertex Vector Search fields <p>The template repo contains sample YAML files for reference</p>	
5	Update the model-config.yaml file	<p>Furthermore, for inference purposes, a separate model-config YAML file captures inference-based configurations. This includes:</p>	

#	Element	Description	Assumptions
	in config/ w/ relevant values	<ul style="list-style-type: none"> • LLM parameters around temperature, top_k, top_p etc • LLM Prompts as YAML fields • ... <p>The template repo contains sample YAML files for reference</p>	
6	Update the templated RAG inference logic	<p>This template repo contains generic code to run a single LLM as part of its logic</p> <p>If a use case requires multiple LLM calls or any additional processing (reranking etc), the same needs to be reflected in the:</p> <ul style="list-style-type: none"> • inference_workflow/src/ folder (contains source scripts for inference & retrieval) • inference_workflow/api/routers/ folder (contains FastAPI logic which imports source scripts & executes them in a predict router) 	
7	Test changes by running locally	<p>Once required changes have been committed, the same can be tested locally by spinning up the Use Case Application locally using <code>make run</code> which spins up the FastAPI server on localhost:8000 by default</p> <p>Local testing can be executed by running a <code>curl</code> command to send API calls directly into the local server</p>	
8	Update unit tests if required	<p>The repo contains unit tests related to all existing functionality in the <code>testing/</code> folder</p> <p>These tests can be updated based on the refactoring as required & tested locally by executing the <code>make pytest</code> command.</p> <p>This Makefile abstracts away execution of key commands & can be referred for other aspects as well</p>	
9	Create a Pull Request	Once relevant changes have been tested & the feature branch is good to go, next logical step would be to create a Pull Request in the repo for review	
10	Merge to main	Post review & approval of the feature branch functionality, the same can be merged into the <code>main</code> branch	

2.5. User Journey #11 - Deploy a RAG Inference Use Case Application & Promote to higher environments

#	Element	Description	Assumptions
	w/ evaluation-based values	<ul style="list-style-type: none"> LLM evaluation on a pre-set Golden Dataset (prompt + expected response) Continuous LLM evaluation based on actual prompt-response pairs generated by the RAG Use Case Application <p>In addition to this generic code to build Vertex pipelines for the above-mentioned 2 patterns, this template abstracts relevant updates to a configuration-based YAML file which exists in the config/ folder</p> <p>This config folder will need a sub-folder w/ naming of the AIB instance (eg vf-uk-aib-prd-cip-ask-lab/nl/live points to vf-uk-aib-prd-cip-ask) & YAML files within said sub-folder pointing to the lab, non-live & live instances</p> <p>For testing purposes, the relevant YAML file needs to be updated w/ evaluation-specific values catering to:</p> <ul style="list-style-type: none"> Project ID Region Evaluation metrics Input Golden Dataset BQ location Expected output BQ location to store evaluation results <p>The template repo contains sample YAML files for reference</p>	pairs already exist in a BQ table as a one-time activity
5	Update code for Vertex pipeline components & the pipeline (if required)	<p>Once the configuration-based YAML files have been updated, base logic in the <code>evaluation/pipeline/</code> folder can be updated if absolutely required.</p> <p>This folder contains generic logic for Vertex Pipeline components that include:</p> <ol style="list-style-type: none"> 1. Extracting data from a BQ table 2. Querying a Use Case app to generate relevant responses 3. Running LLM Evaluation using the Google GenAI Evaluation service 4. Storing evaluation results in the relevant output BQ table <p>For the Golden Dataset Evaluation pattern, all 4 components apply & have been plugged into a generic Vertex pipeline</p>	

#	Element	Description	Assumptions
1	Clone GenAI Template repo	Clone the template repo for your GenAI use case & leverage the boilerplate code & frameworks already provided by the template	
2	Setup a development workspace	<p>Once the assumptions have been met, further setup is required. This includes:</p> <ul style="list-style-type: none"> • Python package installation using poetry as a package manager • Loading environment variables <p>Detailed steps for the same are captured in the CONTRIBUTING.md of the repo</p>	<p>Development workspace has been setup by:</p> <ul style="list-style-type: none"> • Creating a Vertex Workbench instance using this guide - How to : Create a Vertex Notebook - AI Booster • VSCode has been connected to the created Workbench instance using this guide - How to : Connect to AI Booster Notebook vis SSH - AI Booster
3	Create a feature branch	Create a feature branch on the cloned repo to build & test relevant LLM Evaluation - Continuous Evaluation functionality	
4	Update the relevant YAML file/s in config/ w/ evaluation-based values	<p>This template repo contains generic code to execute LLM evaluation catering to 2 patterns:</p> <ul style="list-style-type: none"> • LLM evaluation on a pre-set Golden Dataset (prompt + expected response) • Continuous LLM evaluation based on actual prompt-response pairs generated by the RAG Use Case Application <p>In addition to this generic code to build Vertex pipelines for the above-mentioned 2 patterns, this template abstracts relevant updates to a configuration-based YAML file which exists in the config/ folder</p> <p>This config folder will need a sub-folder w/ naming of the AIB instance (eg vf-uk-aib-prd-cip-ask-lab/nl/live points to vf-uk-aib-prd-cip-ask) & YAML files within said sub-folder pointing to the lab, non-live & live instances</p> <p>For testing purposes, the relevant YAML file needs to be updated w/ evaluation-specific values catering to:</p> <ul style="list-style-type: none"> • Project ID • Region 	<ul style="list-style-type: none"> • Use Case application has been created as part of UJ #10 + #11 & deployed as a Cloud run instance • BigQuery log sink has been created for the Use Case Inference Cloud Run application & contains information on the actual prompt-response pairs generated by the app

#	Element	Description	Assumptions
		<ul style="list-style-type: none"> • Evaluation metrics • Input BQ location containing actual prompt-response pairs generated by a Use Case Application • Expected output BQ location to store evaluation results <p>The template repo contains sample YAML files for reference</p>	
5	Update code for Vertex pipeline components & the pipeline (if required)	<p>Once the configuration-based YAML files have been updated, base logic in the <code>evaluation/pipeline/</code> folder can be updated if absolutely required.</p> <p>This folder contains generic logic for Vertex Pipeline components that include:</p> <ol style="list-style-type: none"> 1. Extracting data from a BQ table 2. Querying a Use Case app to generate relevant responses 3. Running LLM Evaluation using the Google GenAI Evaluation service 4. Storing evaluation results in the relevant output BQ table <p>For the Continuous Evaluation pattern, components 1,3,4 apply & have been plugged into a generic Vertex pipeline</p>	
6	Update unit tests (if required)	<p>The repo contains unit tests related to all existing functionality in the <code>testing/</code> folder</p> <p>These tests can be updated based on the refactoring as required & tested locally by executing the <code>make pytest</code> command.</p> <p>This Makefile abstracts away execution of key commands & can be referred for other aspects as well</p>	
7	Manually compile & submit a Vertex Pipeline job	Once all changes are in place, the Vertex pipeline can be manually submitted via a single python command to test updated functionality	
8a	Validate Vertex Pipeline execution	<p>Once the above-mentioned python command has been executed in a development workspace, the Terminal will provide a link for the running pipeline OR the same can be found in the Vertex Pipeline UI</p> <p>Once submitted & executed, this pipeline should be validated via the UI to verify updated functionality & configuration (if any)</p>	
8b	Validate output BQ table	Once the Vertex pipeline has been submitted & successfully executed, an	

#	Element	Description	Assumptions
		<p>configuration-based YAML file which exists in the config/ folder</p> <p>This config folder will need a sub-folder w/ naming of the AIB instance (eg vf-uk-aib-prd-cip-ask-lab/nl/live points to vf-uk-aib-prd-cip-ask) & YAML files within said sub-folder pointing to the lab, non-live & live instances</p> <p>For promotion purposes, the relevant YAML files for non-live & live need to be updated w/ evaluation-specific values catering to:</p> <ul style="list-style-type: none"> • Project ID • Region • Evaluation metrics • Input BQ location containing actual prompt-response pairs generated by a Use Case Application • Expected output BQ location to store evaluation results <p>The template repo contains sample YAML files for reference</p> <p>Additionally, configuration values setup in the lab YAML as part of pipeline testing in UJ #13 should be used as a reference</p>	
3	Update the cron schedule for the Cloud Scheduler job	<p>In addition to relevant configuration-based YAML files, separate files contain logic to configure the cron schedule of the Cloud Scheduler job to be created/updated</p> <p>cron schedules for both non-live & live Use Case Vendor layer projects will need to be updated separately</p>	
4	Create a Pull Request	Once relevant changes have been tested & the feature branch is good to go, next logical step would be to create a Pull Request in the repo for review	
5a	Merge to main	Post review & approval of the config-update/feature branch functionality, the same can be merged into the <code>main</code> branch	
5b	Create a Release tag	A new release tag needs to be created w/ relevant description supporting it.	

#	Element	Description	Assumptions
		<p>This tag should conform to best practices around tags, specifically of format:</p> <p>vMAJOR.MINOR.PATCH</p> <p>eg - A new release tag can be:</p> <ul style="list-style-type: none"> v1.2.0 	
6a	CI process automatically creates/updates the Cloud Scheduler job in the non-live Use Case Vendor layer project	<p>Once the PR has been pushed to the <code>main</code> branch (see step 5a), the CI/CD pipeline is automatically triggered that:</p> <ul style="list-style-type: none"> Runs relevant unit testing & linting checks Runs the image build & push process & finally outputs an image in Google Cloud Artifact registry of the Use Case Vendor layer build-nl project with a <code>latest</code> tag along with the relevant git commit SHA value Compiles relevant Vertex pipelines & stores them in GCS under a folder w/ name mapping to the git commit SHA value in the build-nl Use Case Vendor layer project Create/update a Cloud Scheduler job for the compiled Vertex pipeline/s in the non-live Use Case Vendor layer project <ul style="list-style-type: none"> This is done via Terraform embedded into the CI/CD pipeline as its final stage & it utilizes the git commit SHA value attached to the image + the compiled pipeline GCS location from the previous stage 	<p>Cloud Build Trigger exists in the Use Case Vendor layer build-nl project to run this CI/CD pipeline</p> <p>This is a one-time manual activity & needs to be implemented once for every GenAI use case</p>
6b	CI process automatically creates/updates the Cloud Scheduler job in the live Use Case Vendor layer project	<p>Once a new release tag has been created (see step 5b), the CI/CD pipeline is automatically triggered that:</p> <ul style="list-style-type: none"> Runs relevant unit testing & linting checks Runs the image build & push process & finally outputs an image in Google Cloud Artifact registry of the Use Case Vendor layer build-live project 	<p>Cloud Build Trigger exists in the Use Case Vendor layer build-live project to run this CI/CD pipeline</p> <p>This is a one-time manual activity & needs to be implemented once for every GenAI use case</p>

#	Element	Description	Assumptions
		<p>with the release tag attached to the image</p> <ul style="list-style-type: none"> Compiles relevant Vertex pipelines & stores them in GCS under a folder w/ name mapping to the release tag in the build-live Use Case Vendor layer project Create/update a Cloud Scheduler job for the compiled Vertex pipeline/s in the live Use Case Vendor layer project <ul style="list-style-type: none"> This is done via Terraform embedded into the CI/CD pipeline as its final stage & it utilizes the release tag attached to the image + the compiled pipeline GCS location from the previous stage 	

1.2.3 User Journeys - Proxy Config Update

- [1. User Journey-specific Background & Context](#)
- [2.1. User Journey #15 - Update Proxy Config for my use case \(Current Approach\)](#)
- [2.2. User Journey #15 - Update Proxy Config for my use case \(Future Roadmap\)](#)

1. User Journey-specific Background & Context

2.1. User Journey #15 - Update Proxy Config for my use case (Current Approach)

