# Tab 1



#### **HOW TO USE**

- (1) Create a copy of this document
- (2) Complete a 'Find and Replace' for the following:

Placeholder	Description
[Customer Legal Name]	Legal customer business name
[Customer Name]	Dba (Doing business as) name for customer
[SOW Date]	Effective Date
[Start Date]	Project start date
[End Date]	Project end date
[Author]	Name and email of Bitstrapped member who wrote SOW

(3) Delete the contents of this page (Page 1)

## Bitstrapped Inc.

[Customer Name]

# Gen Al Model and Model Version Evaluation Framework

[SOW Date]

[Author]



[Customer Logo]

#### **Statement of Work**

This Statement of Work ("SOW") is agreed upon and entered into by and between Bitstrapped Inc. on behalf of itself and its affiliates (hereinafter referred to as ("Bitstrapped") and [Customer Legal Name] (hereinafter referred to as "[Customer Name]" or ("Customer"). This Statement of Work shall be effective on the date of the last signature below, governed by the terms of this SOW, and supersedes any other agreement between the parties related to this subject matter.

## **Contents**

Contents	3
1. Executive Summary	4
1.1. Project Overview	4
1.2. Objectives	5
2. Constraints	6
3. Functional Requirements	7
Non-Functional Requirements	8
4. Scope of Services	9
4.1. Success Criteria	9
4.2. Activities	9
4.3. Deliverables	10
5. Out of Scope	11
6. Roles and responsibilities	12
6.1. Partner Roles and Responsibilities	12
6.2. Customer Roles and Responsibilities	12
7. Dependencies/Assumptions	13
7.1. General assumptions	13
7.2 Customer Responsibilities	13
8. Change Requests	14
9. Project Timeline	15
10. Financials	16
11. Signatures	17

## 1. Executive Summary

This Statement of Work (SOW) covers an implementation pilot of a Generative AI (Gen AI) solution for LLM performance evaluation and observability. The solution enables organizations building Generative AI or LLM-powered products to monitor the performance of their chosen LLM model, evaluate performance across models and model versions, and compare prompt performance optimizations across models.

The objective is to help Customer adopt the Google Cloud Platform (GCP) to take advantage of Google's open approach to Gen Al models and interoperability through Vertex Al Model Garden. This is achieved by streamlining the process of adopting and testing new models and model versions and avoiding issues with model quality or user experience when changing underlying LLM technologies.

The framework implemented as part of this SOW will enable Customer to overcome common challenges with evaluating and selecting optimal large language models (LLMs), such as measuring consistency of results and sensitivity to context. The deployment will include standard testing protocols, automated monitoring, and feedback loops to detect performance drops and support data-driven model decisions.

Built on GCP, the solution will combine native GCP services and open-source tools into a reusable, customizable framework designed to meet observability standards.

### 1.1. Project Overview

Bitstrapped will partner with Customer to deploy the Gen Al evaluation solution, ensuring it is functional in all its components to help address LLM evaluation and performance challenges. Bitstrapped will lead the solution deployment on GCP for a production pilot.

This deployment will include a layered architecture supporting:

- 1. Automated drift detection, with alerting, to flag performance drops.
- Evaluation pipeline to compare the performance of different models. This includes the ability to
  utilize the ground truth or human feedback (if either of these are available) or LLM-based
  evaluation, depending on the availability of data or client-side capture.
- 3. An ingestion point, storage, and integration of user feedback in the evaluation
- 4. Visualization for evaluation insights

The framework will support batch evaluation, giving teams ongoing visibility into model quality and enabling safe, data-driven improvements.

### 1.2. Objectives

- Enhance existing infrastructure to customize a GCP-based Gen Al framework for assessing and monitoring LLM outputs, achieving greater transparency and insight into model behaviour.
- Establish drift detection mechanisms with configurable alert thresholds.
- Create automated evaluation pipelines using both metric-based and LLM-as-judge methodologies.
- Integrate human feedback loops from both domain experts and end-users.
- Provide visualization capabilities for performance metrics and evaluation results.
- Deploy the capabilities using Infrastructure as Code (IaC) with Terraform.

SOW Start Date	[Start Date]
SOW End Date (estimated)	[End Date]

## 2. Constraints

Constraints establish predefined boundaries to manage this engagement's scope, resources, and expectations. The following table outlines the constraints for this engagement:

Constraint	Description	Limit	Inclusion
LLM Use Cases	The # of existing applications or use cases of LLMs supported	1	Deployment will be piloted with one pre-existing LLM solution
Supported Prompts	The # of existing prompts supported	3	Deployment will be piloted for three pre-existing prompts in the supported use case.
Input Data Format	File formats to run evaluation	n/a	CSV or JSONL file(s)
Human Feedback	Format of human feedback supported	Binary or Numeric	An API will provide an ingestion point for human feedback, which will be stored and integrated with evaluation.  It will support binary or numeric feedback scores.  Customer is responsible for integrating the ingestion point to the source of feedback, such as the client side, and/or any migration of baseline data.
LLM Models Evaluated	# of LLM models to evaluate	3	The deployment will be piloted using three suitable LLM models to evaluate performance against
Metrics	# of metrics tracked for evaluation	3	The deployment will produce up to three metrics to be consumed in the visualization
Dashboards	# of visualization dashboards		The deployment will be piloted with 1 dashboard visualization to consume evaluation insights
Observability Platform	# of monitoring systems for tracking model behaviour	1	One of: Arize, Latitude, or Confident-Al, based on requirements analysis  Monitoring setup for 1 to 2 metrics of interest (could be custom or standard metrics, or monitoring the quantitative values representing the human feedback)

## 3. Functional Requirements

This section identifies the corresponding solution components:

ID	Requirement	Solution Component
FR-001	Automated Drift Detection  Monitor drift in the model outputs	Vertex Al Evaluation Service and/or one of the following - Arize, Latitude, Confident-Al
FR-002	LLM Performance Ingestion Ingestion points to integrate existing LLM prompt, response, and model metadata for persistence	TBD
FR-003	Evaluation Pipeline Each LLM model should consume the prompts included to produce the supported metrics for evaluation.  Evaluation will consider: -Drift detection (Batch) -Evaluation (Batch) -Human Feedback (Batch)	Vertex Al Pipelines + GKE + BigQuery
FR-004	LLM-as-Judge Evaluation Secondary LLMs should evaluate model outputs against defined criteria.	GKE + Vertex AI Endpoint
FR-005	Human Feedback Ingestion Ingestion points to ingest, process, and incorporate human feedback.	Cloud Functions + GCS / Firestore / BigQuery
FR-006	Prompt Versioning and Catalogue  Maintain a versioned repository of prompts.	GCS
FR-006	Evaluation History  Record model performance by associating model type, version, prompts, responses, human feedback, and evaluation metrics.	BigQuery
FR-007	Performance Visualization Generate artifacts for reviewing evaluation metrics	BigQuery + Looker or Looker Studio GCS

## **Non-Functional Requirements**

The following section outlines the quality attributes and constraints that the solution will meet:

ID	Requirement	Solution Component	
NFR-001	Scalability – Ability to scale to additional LLM use cases, as well as integrate additional models and sources of human feedback	GKE	
	This should allow Customer to extend this solution to support other LLM use cases.		
NFR-002	<b>Reliability</b> – Ensure consistent evaluation results with error handling and recovery mechanisms for evaluation jobs.	GKE, Cloud Functions	
NFR-003	Data Privacy and Security – Enforce security measures, including data encryption, access control.	Cloud IAM, Cloud Storage Security, Cloud Logging, VPC	
NFR-004	<b>Extensibility</b> – Solution must support Model-agnostic evaluation with minimal code changes.	GKE, Cloud Functions	

## 4. Scope of Services

#### 4.1. Success Criteria

The success criteria define the measurable outcomes and benchmarks that indicate the achievement of project goals and satisfactory delivery of the agreed-upon solution, which include the following:

- Complete a pilot of the LLM evaluation deployment with demonstrated functionality for supported use cases, LLM models, and prompts.
- Automated drift detection produces alerts.
- Evaluation pipelines can be triggered both manually and through event-based automation.
- Access to visualization for comparing and benchmarking the computed performance metrics and evaluation results.
- The deployed solution satisfies the non-functional requirements.
- Documentation for future customization and knowledge transfer will enable the customer team to operate and extend the framework.

### 4.2. Activities

The following section outlines the detailed activities planned to achieve the project objectives and deliverables effectively:

Phase	Activities
Requirements Analysis	<ul> <li>Understand the specific challenges with the Customer LLM deployment that the solution aims to address and the expected outcomes.</li> <li>Identify existing data sources, technology stack, and workflows related to the project scope.</li> <li>Review data sets in scope.</li> <li>Conduct a detailed assessment of evaluation needs and quality standards.</li> <li>Define specific metrics for automated and human-based assessment.</li> <li>Define data schemas for prompt storage, evaluation results, and performance tracking.</li> <li>Design alert thresholds and notification workflow.</li> </ul>
Core Infrastructure Implementation	<ul> <li>Provision required GCP resources and configure access controls.</li> <li>Set up Cloud Storage buckets for prompt cataloging.</li> <li>Configure BigQuery datasets for evaluation metadata and results.</li> <li>Configure base integrations between components.</li> </ul>
Evaluation Pipeline Development	<ul> <li>Deploy a Vertex AI Pipeline for model evaluation.</li> <li>Deploy drift detection and alerting mechanisms.</li> <li>Implement GKE service for LLM-as-judge evaluation.</li> <li>Provision prompt management and versioning system.</li> <li>Build API for feedback ingestion (Batch processing job).</li> <li>Implement visualization of the computed performance metrics and assessment results</li> </ul>

Integration and Testing	<ul> <li>Integrate all components into a cohesive workflow.</li> <li>Integrate existing LLM use case with storage of prompt, response, and metadata</li> <li>Test full evaluation pipelines with sample prompts and models.</li> <li>Validate drift detection with simulated anomalies.</li> <li>Verify monitoring alert propagation and pipeline triggering.</li> <li>Test human feedback integration into the monitoring and evaluation pipelines.</li> </ul>
Production Release	<ul> <li>Environment separation is deployed with Terraform</li> <li>Supported Repository and Git Branches for Prod deployments</li> </ul>
Knowledge Transfer	<ul> <li>Develop comprehensive architecture documentation.</li> <li>Create user guides for framework operation.</li> <li>Prepare customization documentation for future deployments.</li> <li>Conduct weekly knowledge sharing during implementation.</li> <li>Deliver the final solution with a hands-on demonstration of the workflow scenarios.</li> <li>Provide handover documentation with operational instructions.</li> </ul>
Project Management	<ul> <li>Project Planning and Pre-requisites</li> <li>Project Kick-Off meeting</li> <li>Weekly team (internal/external) meetings</li> <li>Weekly project activity and meeting coordination</li> <li>Customer team and stakeholder collaboration</li> <li>Weekly status updates and reporting</li> </ul>

#### 4.3. Deliverables

The customer can expect the following key deliverables from this engagement:

- Full implementation of the evaluation framework with all components as a solution pilot, including:
  - Monitoring implementation with alert configuration.
  - Vertex Al Pipelines for automated model assessment.
  - o LLM-as-Judge service for automated output scoring.
  - o Generation of artifacts for the visualization of metrics.
  - Integrated human feedback
- Fork of Terraform Infrastructure-as-Code deployment
- User guides, operational instructions, and architecture reference.
- Knowledge transfer workshop.
- Results presentation and recommendations

### 5. Out of Scope

The following items are explicitly excluded from the scope of this project:

#### Requirement

#### **Model Training and Retraining**

Training of new foundation models or extensive fine-tuning.

#### **Privacy and Compliance**

The solution must comply with data privacy regulations (e.g., GDPR, CCPA) and ensure the secure handling of user data.

#### **Cloud Foundations**

Set up Google Cloud with foundational best practices for administration.

#### **Production Deployment**

Full-scale deployment and maintenance of the Monitoring solution in a production environment. Deployment will be limited to a pilot in a production environment, for which Customer will be required to manage and support to ensure full operational success.

#### Support and Maintenance

Ongoing support, maintenance, and updates beyond the initial Proof of Concept (PoC) phase.

#### HITI Interfaces

Custom development of sophisticated HITL interfaces.

#### Integration

Integration with existing customer monitoring systems, non-GCP environments or cloud platforms.

#### **Custom Development**

Custom development of complex evaluation metrics is not specified in the requirements.

#### **BI Foundations**

Any setup, foundations, or integrations required for Looker or Looker Studio

#### Third-party tool integration

Integration with a non-GCP tool is out of scope

## 6. Roles and responsibilities

## **6.1.** Partner Roles and Responsibilities

Role	Role Description
Project Manager	<ul> <li>Serve as the primary point of contact for the Customer</li> <li>Develop and maintain project plans, timelines, and status reports</li> <li>Coordinate project activities and deliverables across all teams</li> <li>Facilitate regular project status meetings with the Customer</li> <li>Manage project risks, issues, and changes in scope</li> <li>Ensure project deliverables meet the Customer requirements</li> </ul>
ML Engineer	<ul> <li>Review provided dataset consisting of prompts and responses.</li> <li>Design data schemas for evaluation metrics and implement BigQuery integration.</li> <li>Configure visualization queries</li> <li>Lead LLM integration and evaluation, and implement model serving components</li> <li>Design evaluation methodologies</li> <li>Configure LLM-as-judge functionality</li> </ul>
Cloud Engineer	<ul> <li>API implementation</li> <li>Supporting infrastructure</li> <li>Build cloud infrastructure</li> <li>Deploy GCP resources</li> <li>Configure Vertex AI Pipeline components</li> <li>Implement pub/sub and event triggers</li> <li>Deploy GKE evaluation services</li> <li>Set up the observability tools integration</li> </ul>

## **6.2.** Customer Roles and Responsibilities

Customer Role	Responsibilities
Project Sponsor	<ul> <li>Serve as the end decision-maker</li> <li>Act as a point of escalation for issues</li> <li>Sign off on deliverables</li> <li>Provide business criteria for data exploration</li> <li>Identify priority predictive outcomes</li> </ul>
Project Manager	<ul> <li>Align Customer resources to support successful project completion</li> <li>Help resolve any issues or questions in a timely manner</li> </ul>
SME such as Product Manager or Technical Lead	<ul> <li>Collaborate with Bitstrapped to optimize Gen AI for the use case</li> <li>Collaborate with Bitstrapped on vector database configuration</li> </ul>

## 7. Dependencies/Assumptions

## 7.1. General assumptions

- The project timeline was based on information interpreted from related conversations with Customer to date. To the extent material discrepancies are uncovered, this SOW may need to be altered to accommodate timeline, resources, and/or cost changes.
- Customer will provide mutually agreed-upon personnel (IT, etc.) and resources (documentation, systems, etc.) on time.
- Customer Project sponsors will help resolve any issues or questions related to the work on time.
- The Customer team will assist with acquiring and implementing the Google Cloud environment and other tools required for development and deployment.
- Google Cloud usage fees are not included in the costs specified in the Financials Section. Customer shall be responsible for all Google Cloud usage fees.
- The new system has no compliance / regulatory / data location requirements for customer data.
- Partner does not warrant or otherwise represent the functionality of any vendor product used during this project. The respective vendor must provide such representations.
- Work related to this SOW is expected to be 100% remote. Travel to customer sites and on-site services are not required.
- This SOW is based upon the accuracy of the information the Customer provides.
- Any additions, changes, or modifications to the Services specified herein are outside the scope of this SOW.
- All Partner work will be performed remotely at Partner's facilities or resource locations.
- Customer shall make access for Partner resources available to the systems, applications, and knowledgeable personnel during designated time frames, which will be established during the project kick-off meeting. Failure to provide this timely access could delay the completion of the services.
- Customer already has Gen AI models serving inference and has the ability to add additional models as part of this evaluation framework deployment, and/or is using Bitstrapped to do this work as part of a separate SOW

### 7.2 Customer Responsibilities

- Align stakeholders to attend the kickoff meeting and regular weekly calls as part of the engagement; provide timely feedback on deliverables.
- Identify and ensure the availability of subject matter experts with the necessary skills and information regarding the technical infrastructure and the technical and functional application technologies.
- Customer shall be responsible for addressing any internal dependencies (including, but not limited to, business approval, network and security approval)
- Customer shall assign a named Project Sponsor who will be working with Partner.
- Identify data and content conducive to the knowledge base and model goals
- Provide knowledge sharing on datasets, including domain expertise and pre-processing requirements.
- Provide business and end-user requirements for Gen Al.
- Participate in training, feedback, and optimization exercises.
- Participate in regular project status meetings and engage in knowledge transfer sessions and

workshops to understand the end-to-end solution.

- Assume responsibility for the ongoing maintenance and management of the solution post-deployment.
- Deploy any Gen Al models to serve inference
- Maintain a prompt catalogue and versions
- Data preparation and cleaning
- Data integrations or third-party tool integrations
- Scaling out evaluation pipelines, observability dashboards, and other functional requirements beyond the limits of the Constraints section
- Integration with a human feedback source

### 8. Change Requests

Any change to this SOW shall not take effect unless and until a Change Request to this SOW is fully executed by Customer and Partner. The Change Request shall include the following information:

Date of Master Services Agreement (Optional)	
Date of Change Request	
Name and Date of Impacted SOW	
Description of modified Services	
Description of modified Deliverables	
Impact on resources	
Impact on Project Timeline	
Resulting change in Cost	
Effective Date of approval	

## 9. Project Timeline

The project timeline is based on an estimated start date of approximately two weeks after all parties sign this SOW. The start date is subject to change based on the availability of critical team members from the Customer and Partner.

The anticipated project duration is approximately **8 weeks**. This estimate assumes timely reviews and/or approvals from Customer.

Deliverable	WK 01	WK 02	WK 03	WK 04	WK 05	WK 06	WK 07	WK 08
Requirements Analysis and Customization								
Core Infrastructure Implementation								
Evaluation Pipeline Development								
Integration and Testing								
Pilot Release								
Documentation & Knowledge Transfer								

### 10. Financials

Partner will provide the Services on a fixed fee basis per the table below, including the total DAF amount as a line item.

Service	Cost
Requirements Analysis	\$3,248
Core Infrastructure Implementation	\$9,744
Evaluation Pipeline Development	\$19,488
Integration and Testing	\$9,744
Release and Deployment	\$9,744
Documentation & Knowledge Transfer	\$3,248
Project Management	\$3,248
Subtotal	\$58,464
Partner Services Funds (PSF)	\$58,464
Total Due from Customer	\$0

**Google DAF Investment.** Customer acknowledges and agrees that the Google DAF Investment is contingent upon meeting the Success Criteria shown in this document. The total of all Google payments shall not exceed the predetermined amount for the selected Consumption Pack. If Customer terminates this SOW or does not permit Partner to meet the Success Criteria for any reason, then Customer will be obligated to pay Partner under this SOW.

**Estimated Project Expenses.** Expenses, if any, would be billed to Customer as incurred with prior approval from Customer and per Customer policies. No travel or travel-related expenses shall apply to services under this SOW.

## 11. Signatures

The parties have caused this SOW to be executed by their duly authorized representatives as of the effective date set forth above.

[Customer Name]	Bitstrapped Inc.
Signature:	Signature:
Name:	Name:
Title:	Title:
Date:	Date: