# Enhancement & Customizations for

## Workshop: RAG using Amazon Bedrock Agents and Knowledge Base

## Description:

The original "**RAG using Amazon Bedrock Agents and Knowledge base**" workshop guides you how to create a GenAI bedrock agent that can answer questions about Amazon SageMaker.

This document will guide you how to enhance this workshop by:

1) Use your own data as the data source for the Bedrock agent.
2) Customize your own Agent & Knowledge Base instructions (Prompts)
3) Creating a webpage that will allow users from outside of this workshop environment to experiment with your deliverables and send prompts to the agent that you created.
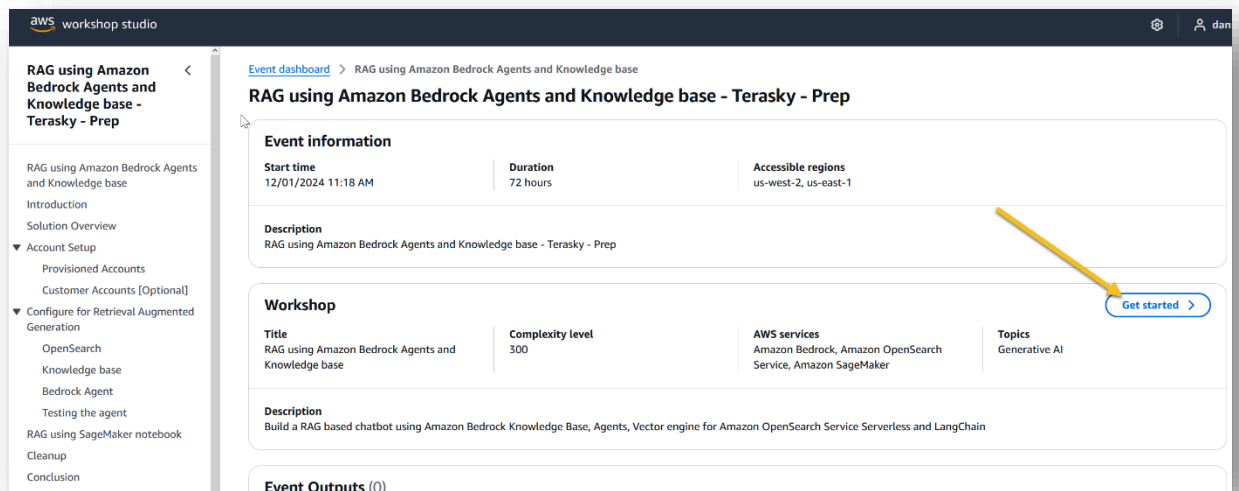
# Use your own data

The workshop session starts with the deployment of the "**rag-w-bedrock-kb**" CloudFormation stack. (In AWS hosted events this stack is pre-deployed before you login to the workshop account)

The **rag-w-bedrock-kb** stack is used to define resources like… IAM roles, Access Policies, S3-Bucket and more… that you get to use in the workshop tasks.

The **rag-w-bedrock-kb** stack also import SageMaker-documentation HTML files into the S3-Bucket, and these files are used as the data source for creating a vector database that the bedrock agent will use.
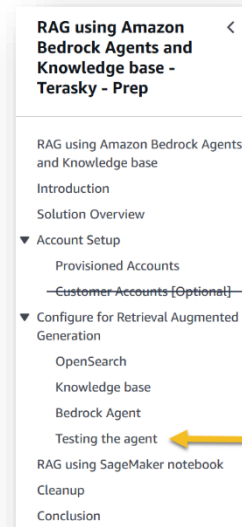
In order to use your own data in this workshop, we need to delete these SageMaker-documentation files and upload your files to the same S3-Bucket.

1) Once you have uploaded the relevant files to the S3-Bucket. Go to the Event dashboard page and click - **Get started >**



2) Follow the workshop lab instructions till you complete the **Testing the agent** phase.

   **Customer Accounts [Optional]** phase is not required.



   **Important:**
   In these workshop steps you will be asked to define resources with specific names like "**sagemaker-qa**" (Bedrock agent)  or "**sagemaker-docs**" (Bedrock Knowledge Base).

   Please make sure that you use the exact names because the initial permissions rules that were defined in this environment are using these names.

# Define your own agent & knowledge Base instructions

On step **Bedrock Agent** tasks 2 & 4 – you are asked to define
**Instructions for the Agent** & **Knowledge base instructions for agent**

The suggested prompts are for a Q&A agent that answer to SageMaker related questions.

This is the time to get creative –

- What prompt will you use?
- What LLM works best for you?
- How can you make your prompt helpful for users?
- Will the prompt hold any output examples?

**Troubleshooting: If your prompt is not working like it should…**

- Are you getting this info message?

> (i) To optimize agent performance, the agent doesn't use instructions if all of the following conditions are true:
>
> - This agent has no action groups
> - There is only one knowledge base associated to this agent
> - No advanced prompts are overridden
> - User input and code interpreter are disabled
>
> If any of the above conditions is false, the agent will use the instructions.Learn more

In order to overcome this limitation let's add a second knowledge base to our agent.

- o On the bedrock page, open the Knowledge base pane
- o Click **Create Knowledge base**



- o Name the knowledge base "knowledge-empty" and click **Next**

**Provide knowledge base details**

**Knowledge base details**
Knowledge base name
knowledge-empty
Valid characters are a-z, A-Z, 0-9, _ (underscore) and - (hyphen). The name can have up to 50 characters.
Knowledge base description – *optional*
Enter description
Valid characters are a-z, A-Z, 0-9, _ (underscore) and - (hyphen). The name can have up to 200 characters.

o   Name the Data source "data-source-empty", Select the S3-Bucket and click **Next**



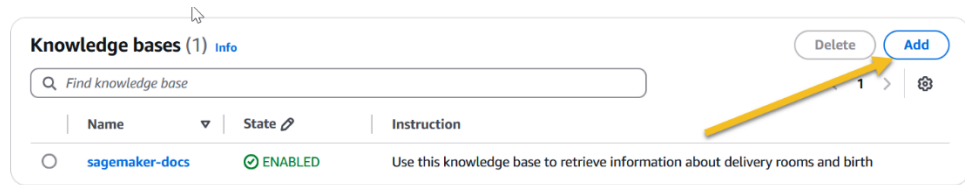o   Select the **Titan Text Embeddings V2** model and click **Next**.
(Without setting the vector DB)



o   Click **Create Knowledge base** and wait for the creation to be completed
o   Navigate back to the **Bedrock → Agents** pane
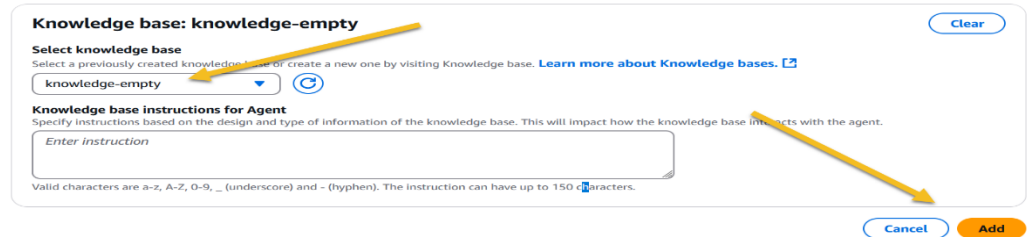o   Select the **sagemaker-qa** agent and click **Edit**

o Under the **Knowledge base** pane click **add**



o Select the empty Knowledge base that you created and click **Add**



o Click **Save and Exit** in the Agent Builder page
o Click **Prepare**
o Test your agent…
o Still not getting the reply that you are expecting…
  ▪ Click **Edit in Agent Builder**



  ▪ Click **change** and select a different model



  ▪ **Save** the agent, **Prepare** the agent and test it
  ▪ Experiment with different models till you get the required results

<u>Prompt example:</u>

- **Instructions for the Agent:**
  You are a medical assistant by the name Daniela.
  Your task is to answer patient's questions based on the data in the knowledge base only.

  The answer should be written in simple language in an E-mail format with new lines between sentences to format a well-structured email.

  The context of the questions is always around birth and delivery rooms.

  Example email reply:
  Dear patient,

  It is recommended that the expectant mother eat snacks that she enjoys during childbirth,
  as the labor process can be long and the mother may feel hungry.
  However, heavy meals are not recommended, especially if the mother is receiving an epidural,
  as she can only drink water during that time.

  I recommend packing light snacks in the delivery room bag to help sustain the mother's energy during labor.

  Thank you and see you at the delivery room,
  Daniela

- **Knowledge base instructions for agent:**
  Your task is to generate email based on the patient's query. Answers questions that are only based on the knowledge base and rephrase the answer in simple English.

# Congratulations!

You have created a GenAI bedrock agent with an OpenSearch vector Database that hosts your embedded vectorized data as a RAG.
Now…
Let's expose your agent to the world

# Expose your Agent to the world

This section will instruct you on how to deploy a CloudFormation template that will create a webpage that will allow external users to send questions to & get-replies from your Bedrock Agent from outside the workshop environment.

In order to access our Bedrock agent with a lambda function, we have to define an Agent Alias.

1) Navigate to the **Bedrock → Agents** pane
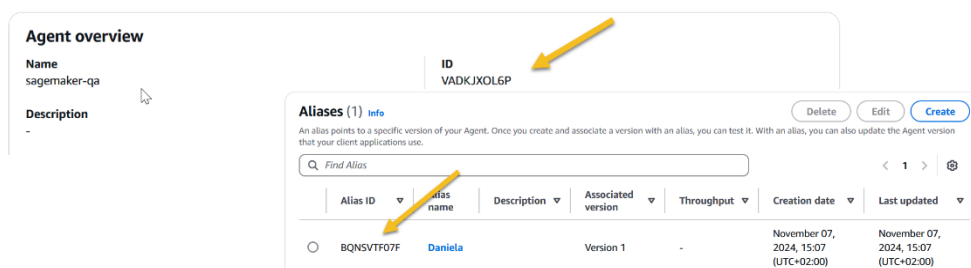2) Click on the **sagemaker-qa** agent



3) Scroll down to the Aliases pane and click **Create**



4) Type an **Alias name** and click **Create Alias**



5) From the agent page copy the **Alias ID** & the **Agent ID**

**Deploy the Agent-Frontend CloudFormation**

1) Navigate to the CloudFormation page and click **Create Stack** → **With new resources**



2) Select **Choose an existing template** & **upload a template file**, and click Choose file



3) Select the stack file that you got from the workshop facilitator and click open



4) Click **Next**

5) Fill the required information
   a. **Stack name**:        Agent-Frontend
   b. **AgentAliasId**:       <The bedrock agent alias ID that we saved before>
   c. **AgentId**:            <The bedrock agent ID that we saved before>
   d. **ProductDescription**: <A description of the product>
   e. **ProductName**:        <The name of the product>
   f. **TeamName**:           <Your Name/s>

   And click **Next**



6) Acknowledge and click **Next**



7) Review and click **Submit**

8) Wait 2-3 minutes till all the resources get created and navigate to the Lambda page

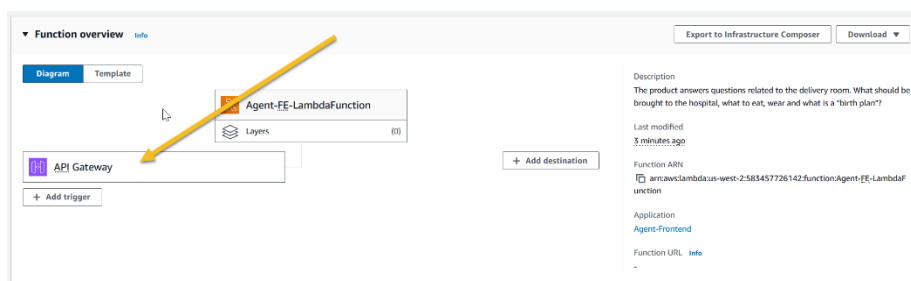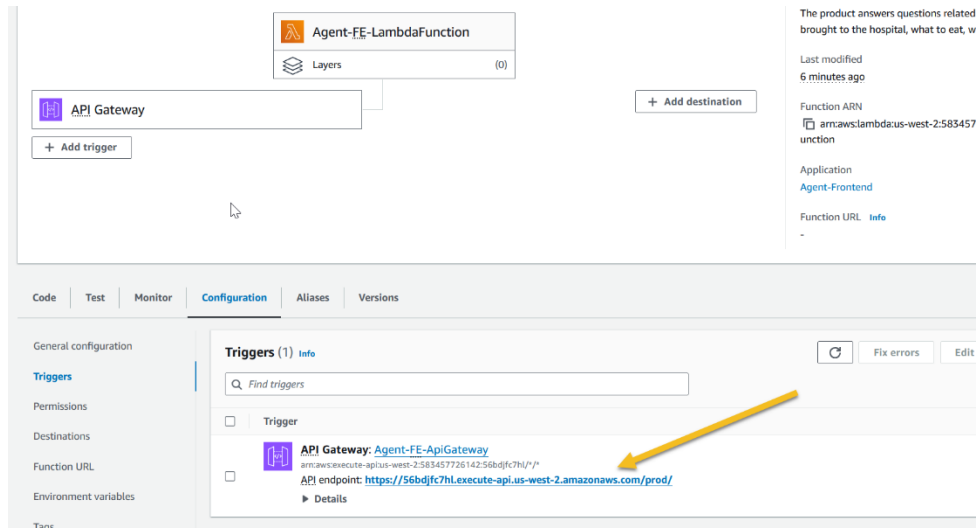9) In the Lambda page, click on the **Agent-FE-LambdaFunction** function
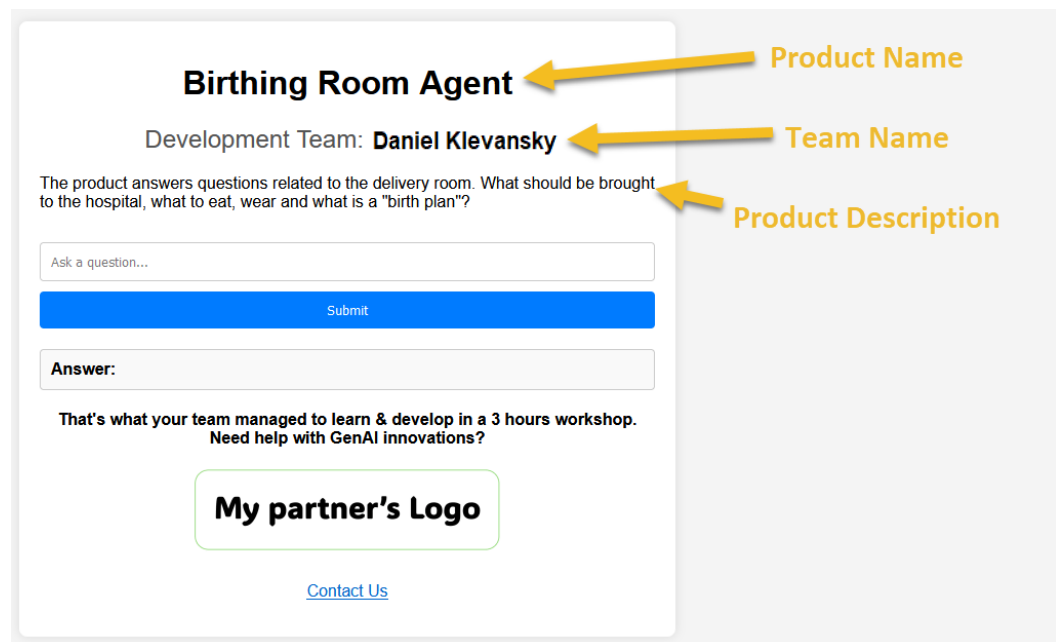


10) Click on the API Gateway component

11) Under configuration click on the API endpoint link and access your agent via the internet.



12) You should see a web page that looks like this…



13) Access this link from your phone & enjoy your first Bedrock Agent