# Practical - 2

**Aim :-** Implementation of a RAG (Retrieval-Augmented Generation) using AWS

## STEPS:

### Step 1: Setup AWS Services

1. **Create an AWS Account**: Sign up for AWS using a debit/credit card, email, and phone number.
2. **Enable Amazon Bedrock**: Navigate to AWS Bedrock and request access.
3. **Select LLM for Generation**: Choose the Meta model (Llama 3 8B Instruct) for text generation.
4. **Enable Amazon Titan Embeddings**: Use Titan Embeddings to generate vector representations of documents.
5. **Set Up Amazon S3**: Store raw text documents and datasets in an S3 bucket.

### Step 2: Data Preprocessing & Embeddings

6. **Upload Documents to S3**: Store knowledge base files (PDFs, TXT, etc.) in an S3 bucket.
7. **Generate Embeddings**: Use Amazon Titan Embeddings to convert documents into vector representations.
8. **Store Embeddings in Amazon OpenSearch**: Use OpenSearch for vector-based similarity search.

### Step 3: Retrieval System & Model Integration

9. **Implement a Search Pipeline**: Use OpenSearch's kNN search to retrieve relevant documents.
10. **Context-based Generation**: Pass retrieved results to Amazon Bedrock's Meta model to generate responses.

### Step 4: API & Frontend Deployment

11. **Develop an API**: Use AWS Lambda & API Gateway to create a chatbot backend.
12. **Deploy the Chatbot Frontend**: Use AWS Amplify or Amazon S3 with CloudFront to host the chatbot UI.

## Step 5: AWS Configuration & Deployment

13. **Configure AWS CLI**: Run the following command and enter your AWS credentials:

aws configure

- Enter **Access Key ID** and **Secret Access Key**.
- Set default **region** (e.g., us-west-2).
- Select output format as json.

14. **Test the Chatbot**: Deploy the API and frontend, then run queries to test retrieval-augmented responses.

# Identity and Access Management (IAM)

Q Search IAM

**Access management**
- Dashboard
- User groups
- Users
- Roles
- Policies
- Identity providers
- Account settings
- Root access management  New

**Access reports**
- Access Analyzer
- External access
- Unused access
- Analyzer settings
- Credential report
- Organization activity

CloudShell    Feedback

## Users (1) Info

An IAM user is an identity with long-term credentials that is used to interact with AWS in an account.

Q Search

| User name | Path | Groups | Last activity | MFA | Password age | Console last sign-in |
|---|---|---|---|---|---|---|
| premansh-user1 | / | 0 | 23 minutes ago | - | Yesterday | February 06, 2025, 19:... |

C    Delete    Create user

< 1 >

aws ::: Q Search [Alt+S] 🔍 ⊡ 🔔 ? ⚙ 🌐 United States (Oregon) ▾ | stargem ▾

🔵 Amazon S3

**Amazon S3** ✕ ‹

General purpose buckets
Directory buckets
Table buckets
Access Grants
Access Points
Object Lambda Access Points
Multi-Region Access Points
Batch Operations
IAM Access Analyzer for S3

Block Public Access settings for this account

▾ **Storage Lens**

Dashboards
Storage Lens groups
AWS Organizations settings

Feature spotlight 🔢

▶ **Account snapshot** - *updated every 24 hours*  `All AWS Regions`

Storage Lens provides visibility into storage usage and activity trends. Metrics don't include directory buckets. *Learn more* ↗

**View Storage Lens dashboard**

General purpose buckets | Directory buckets

**General purpose buckets** (2) Info  `All AWS Regions`

Buckets are containers for data stored in S3.

🔄   ⎘ Copy ARN   Empty   Delete   **Create bucket**

‹ 1 › ⚙ ▾

| Q Find buckets by name | | |
|---|---|---|
| ▲ Name ▼ | AWS Region ▼ | Creation date |
| ⚪ myragknowledgebase | Europe (Stockholm) eu-north-1 | February 5, 2025, 18:51:17 (UTC+05:30) |
| ⚪ premanshragbucket | US West (Oregon) us-west-2 | February 6, 2025, 18:23:08 (UTC+05:30) |

| IAM Access Analyzer |
|---|
| View analyzer for eu-north-1 |
| View analyzer for us-west-2 |

🔲 CloudShell  Feedback

United States (Oregon) ▼  ☐  stargem ▼

Amazon S3 > Buckets > premanshragbucket

# Amazon S3

**General purpose buckets**
Directory buckets
Table buckets
Access Grants
Access Points
Object Lambda Access Points
Multi-Region Access Points
Batch Operations
IAM Access Analyzer for S3

Block Public Access settings for this account

▼ **Storage Lens**
Dashboards
Storage Lens groups
AWS Organizations settings

Feature spotlight  11

## premanshragbucket Info

| Objects | Metadata | Properties | Permissions | Metrics | Management | Access Points |

### Objects (2)  ⟳

Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory ↗ to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. Learn more ↗

☐ Copy S3 URI    ☐ Copy URL    ↓ Download    Open ↗    Delete    Actions ▼    Create folder    ↑ Upload

🔍 Find objects by prefix

◯ Show versions

< 1 > ⚙

| ☐ | Name ▲ | Type ▼ | Last modified ▼ | Size ▼ | Storage class ▼ |
|---|--------|--------|-----------------|--------|-----------------|
| ☐ | 📄 azw(1).pdf | pdf | February 6, 2025, 18:25:16 (UTC+05:30) | 4.4 MB | Standard |
| ☐ | 📄 azw(18).pdf | pdf | February 6, 2025, 18:25:21 (UTC+05:30) | 1.1 MB | Standard |

# Amazon Bedrock

▼ **Getting started**
  Overview
  Providers

▼ **Foundation models**
  Model catalog  New
  Marketplace deployments  New
  Custom models (fine-tuning, dist...)
  Imported models
  Prompt Routers  Preview

▼ **Playgrounds**
  Chat / Text
  Image / Video

▼ **Builder tools**
  Agents
  Flows
  Knowledge Bases
  Prompt Management

▼ **Safeguards**
  Guardrails

## Create a Knowledge Base

- **Knowledge Base with vector store:** Build a fully customizable Knowledge Base with maximum flexibility. Specify the location of your data, select an embedding model, and configure a vector store. Bedrock stores and updates your embeddings.
- **Knowledge Base with structured data store:** Build a Knowledge Base which can connect to a structured data source.
- **Knowledge Base with Kendra GenAI Index** - *new:* Build a Knowledge Base powered by Kendra GenAI Index, offering out-of-the-box high semantic accuracy and the flexibility to reuse the index across Amazon Q Business and Amazon Bedrock Knowledge Bases.

Create ▼

## Test the Knowledge Base

Query your Knowledge Base in the test window. You can get source text chunks, or you can use the chunks to get responses from a foundation model.

## Use the Knowledge Base

Integrate your Knowledge Base into your application as is or add it to agents.

## Knowledge Bases (1)

Edit   Delete   Test Knowledge Base   Evaluate   Create ▼

Q Find Knowledge Base

< 1 > ⚙

| | Name ▽ | Status | Type ▽ | Data so... ▽ | Source ... ▽ | Descrip... ▽ | Creatio... ▽ | Last sy... ▽ | Last sync ▽ |
|---|---|---|---|---|---|---|---|---|---|
| ○ | premansh... | ✓ Available | Vector store | 1 | 2 | - | February ... | - | February ... |

stargem ▼

United States (Oregon) ▼

**Amazon Bedrock**

Amazon Bedrock > Knowledge Bases > premansh-knowledgebase1

# premansh-knowledgebase1

Test | Delete

▼ **Getting started**

Overview

Providers

▼ **Foundation models**

Model catalog  New

Marketplace deployments  New

Custom models (fine-tuning, dist...)

Imported models

Prompt Routers  Preview

▼ **Playgrounds**

Chat / Text

Image / Video

▼ **Builder tools**

Agents

Flows

Knowledge Bases

Prompt Management

▼ Safeguards

## Knowledge Base overview

Edit

**Knowledge Base name**
premansh-knowledgebase1

**Knowledge Base description**
—

**Service Role**
AmazonBedrockExecutionRoleForKnowledge
Base_uop9r [↗]

**Log Deliveries**
Configure log deliveries and event logs in the
Edit page.

**Retrieval-Augmented Generation (RAG)
type**
Vector store

**Knowledge Base ID**
[🗐] AIS2LJY2OP

**Status**
⊘ Available

**Created date**
February 06, 2025, 19:14 (UTC+05:30)

## Data source (1)

Sync | ⊖ Stop sync | Add | ▶

Data sources contain information returned when querying a Knowledge Base.

strategy for your Knowledge
Base, select the configurations
icon ⚏ .

◆◆  what is aws?

🔆  AWS, or Amazon Web Services, is
a cloud computing platform that
allows users to rent virtual
servers, storage, databases, and
other services over the internet. It
provides a scalable and flexible
way to build and run applications,
with a wide range of services that
can be used to create complete
solutions. [1][2]
Show details >

△ Run

*Enter your message here*

⚡ CloudShell   Feedback