



Accelerate Gen AI workloads with AWS Serverless Compute

Build Gen AI Applications on serverless architectures

Dr. Chom Trevai

Generative AI Business Development -
ASEAN

Agenda

Gen AI Ecosystem

Why build Gen AI on AWS Serverless compute

Use cases for Gen AI with Serverless

What is Generative AI?



AI that can
generate content
close enough to human created
content for real-world tasks



Powered by
foundation models
pre-trained on large sets of data with
several hundred billion parameters



Applicable to
many use cases
like text summarization, question
answering, digital art creation,
code generation, etc.



Tasks can be
**customized for
specific domains**
with minimal fine-tuning



New Volvo car concept design by midjourney
Credit: @sugardesign_1 Instagram

Generative AI Stack

APPLICATIONS THAT LEVERAGE LLMs AND FMs



Amazon Q



AWS App Studio

TOOLS TO BUILD WITH LLMs AND OTHER FMs



Amazon Bedrock

Guardrails | Agents | Studio | Customization | Custom Model Import | Amazon Models

INFRASTRUCTURE FOR FM TRAINING AND INFERENCE



GPUs



Trainium



Inferentia



SageMaker



UltraClusters



EFA



EC2 Capacity Blocks



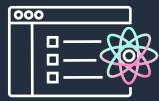
Nitro



Neuron



Generative AI can be used for a wide range of use cases



Enhance customer experience

CHATBOTS

VIRTUAL ASSISTANTS

AI-POWERED CONTACT CENTER

PERSONALIZATION



Boost employee productivity

CONVERSATIONAL SEARCH

SUMMARIZATION

CODE GENERATION

DATA TO INSIGHTS



Creativity & content creation

WRITING

MEDIA

DESIGN

MODELING



Improve business operations

DOCUMENT PROCESSING

PROCESS OPTIMIZATION

CYBERSECURITY

DATA AUGMENTATION

Gen AI ecosystem

Target Candidates for Gen AI with Serverless

Consumers

Use existing pre-trained Foundation models. Typically via API.

Tuners

Customize pre-trained Foundation models. Typically SaaS and customizing per customer. Focus on inference

Builders/ Providers

Pre-train their own foundation models from scratch. Invest heavily in compute (Stability AI, Bloom, AI21 Labs, ...)

Building generative AI applications is challenging



Accessing
multiple FMs
and newer
versions



Customizing
FMs is not easy



Data privacy
and security



Getting FMs
to execute tasks



Connecting to
data sources



Difficult
to manage
infrastructure

What does the future hold for generative AI?

Agents

Multi-modal

Multiple models

AI Policies & Standards



Why build Gen AI on AWS Serverless Compute?

Combining Speed and Power

Speed of Serverless

Power of Gen AI

Rapid delivery
of smarter
applications
and features
with focus on
Innovation



AWS
Lambda



Amazon
ECS



AWS
Fargate



AWS
Step
Functions



Amazon
EventBridge



Amazon
SageMaker



Amazon
Bedrock



Amazon Q for
Developer



Consuming Foundation models - Chatbots and Virtual Assistants

Model invocation - Synchronous



FM Endpoint

- Amazon Bedrock
- Amazon SageMaker

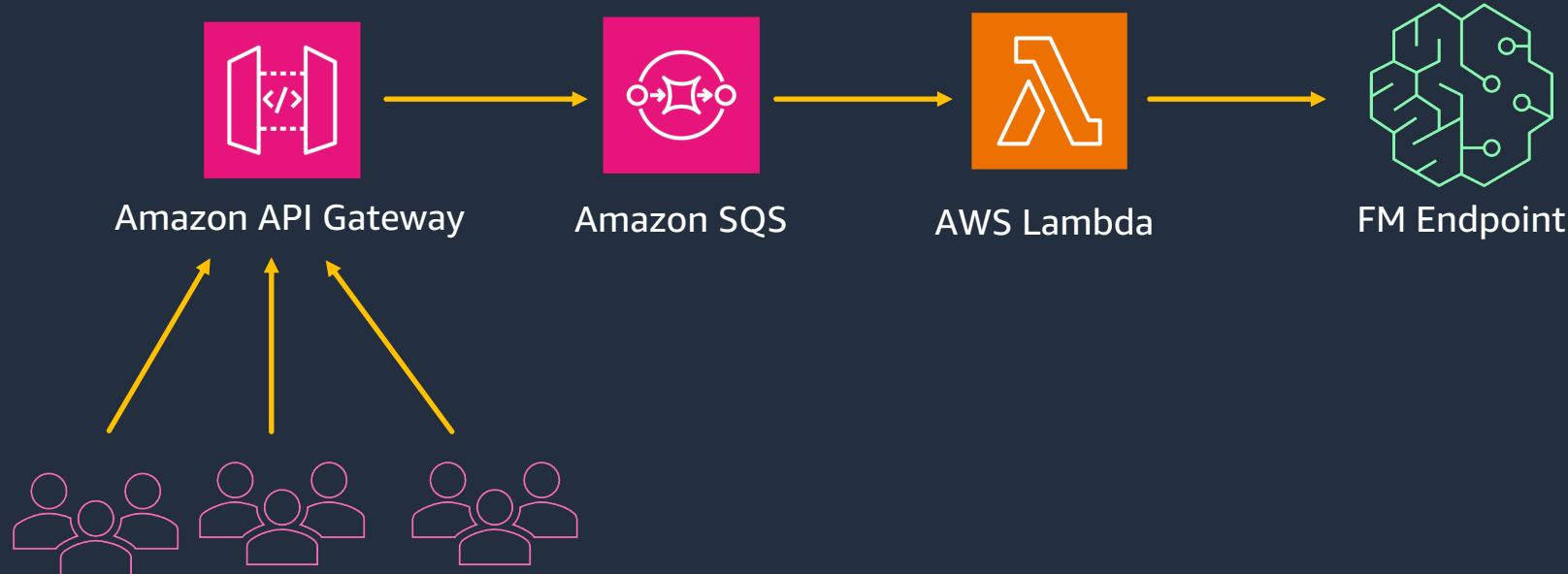
Serverless Services

- Auth, Rate limiting, Caching with Amazon API Gateway
- AWS Lambda calling an endpoint

Use cases

- Simple Q/A
- Content generation
- Text Summarization

Model invocation at scale - Asynchronous



FM Endpoint

- Amazon Bedrock
- Amazon SageMaker

Serverless Services

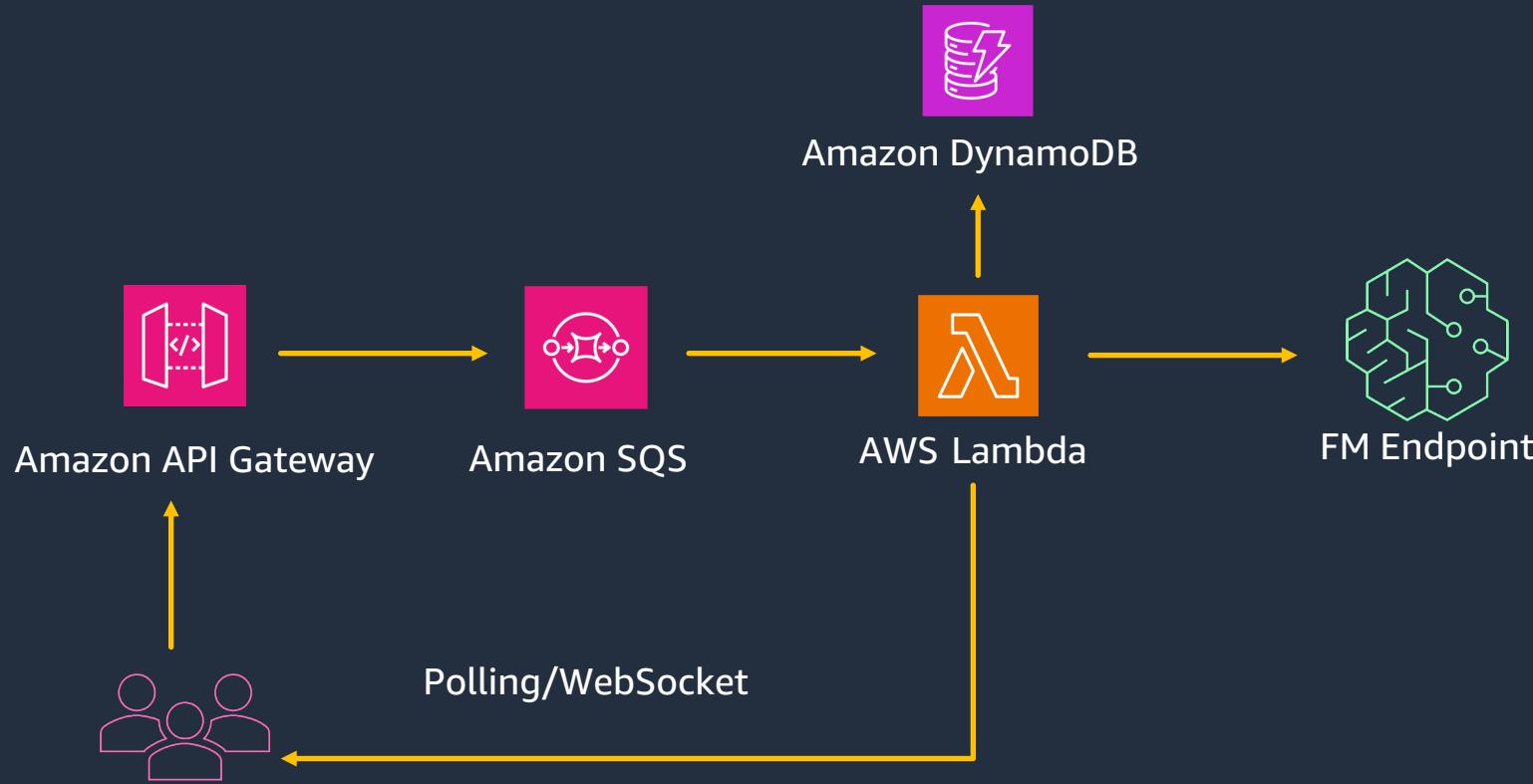
- Auth, Rate limiting, Caching with Amazon API Gateway
- Amazon SQS queue provides control over the scaling
- AWS Lambda calling an endpoint

Use cases

- Simple Q/A
- Content generation
- Text Summarization

Conversation History

Model invocation with conversation memory



FM Endpoint

- Amazon Bedrock
- Amazon SageMaker

Serverless Services

- Auth, Rate limiting, Caching with Amazon API Gateway
- Amazon SQS queue provides control over the scaling
- AWS Lambda calling an endpoint
- Use polling, WebSocket, or IoT core topic for response
- Store conversation history
- Co-ordinate multiple tasks with less code

Use cases

- Virtual Assistants
- Text Summarization
- Context based Q/A

Build context based on Internal knowledge

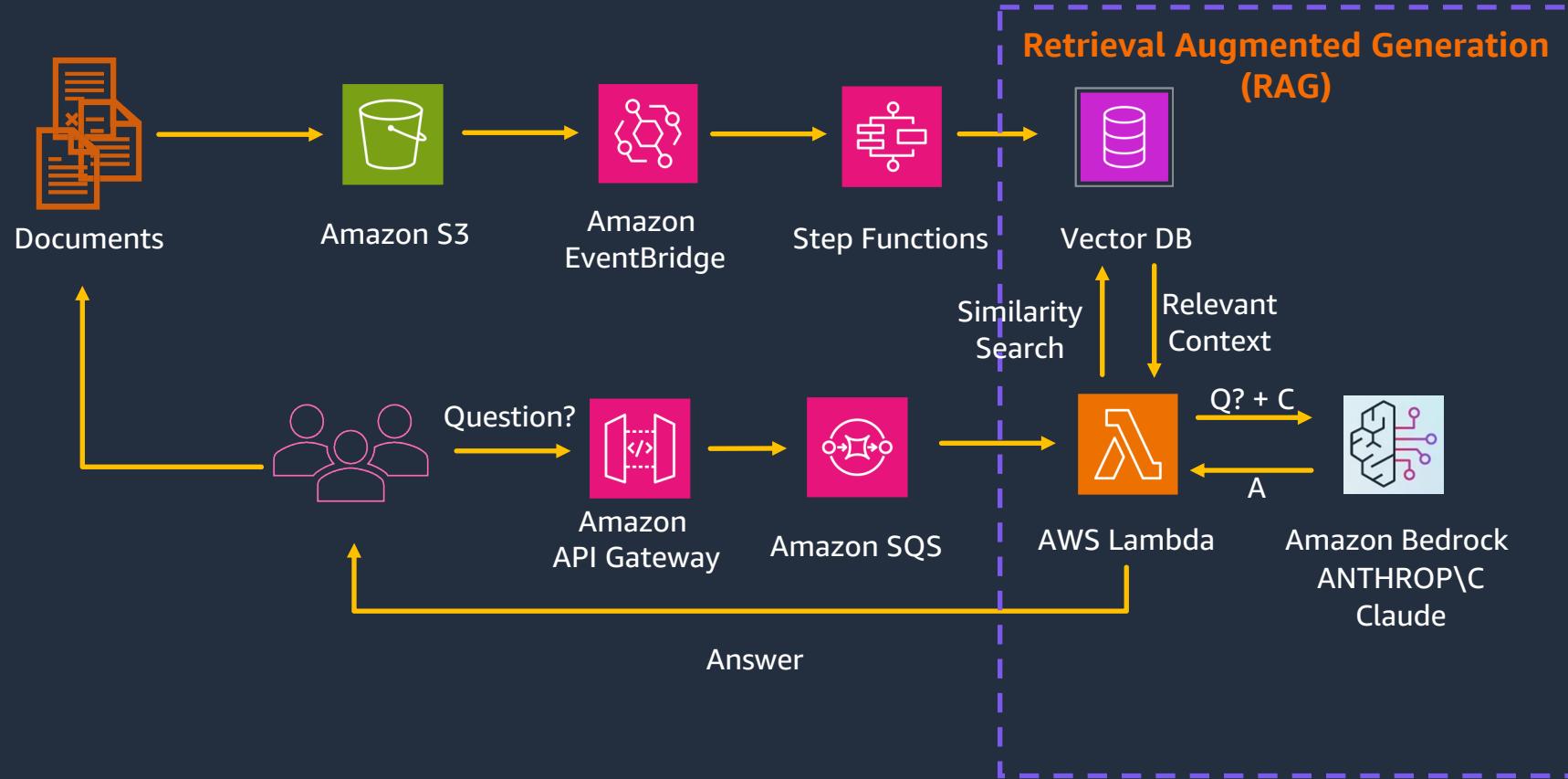


Vector embeddings support RAG



LLM invocation with RAG

Retrieval Augmented Generation (RAG)



FM Endpoints

- Amazon Bedrock
- Amazon SageMaker

Vector Databases

- Amazon OpenSearch
- Chroma DB
- Pinecone

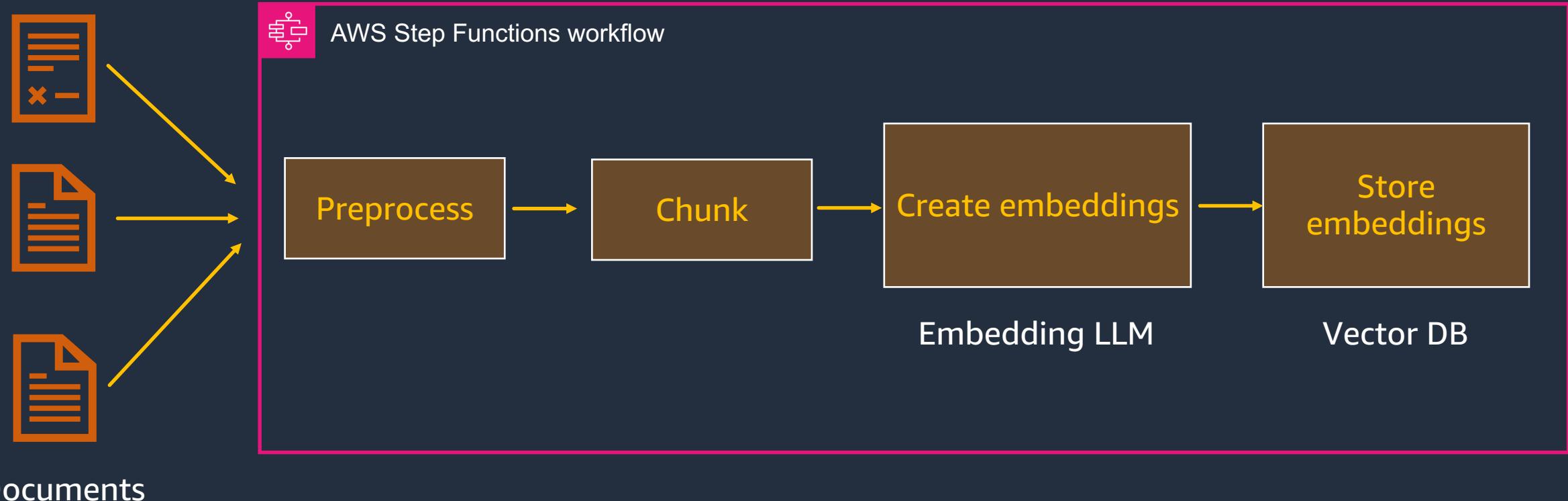
Serverless Services

- Native integration with Amazon S3
- Orchestration with Step Functions

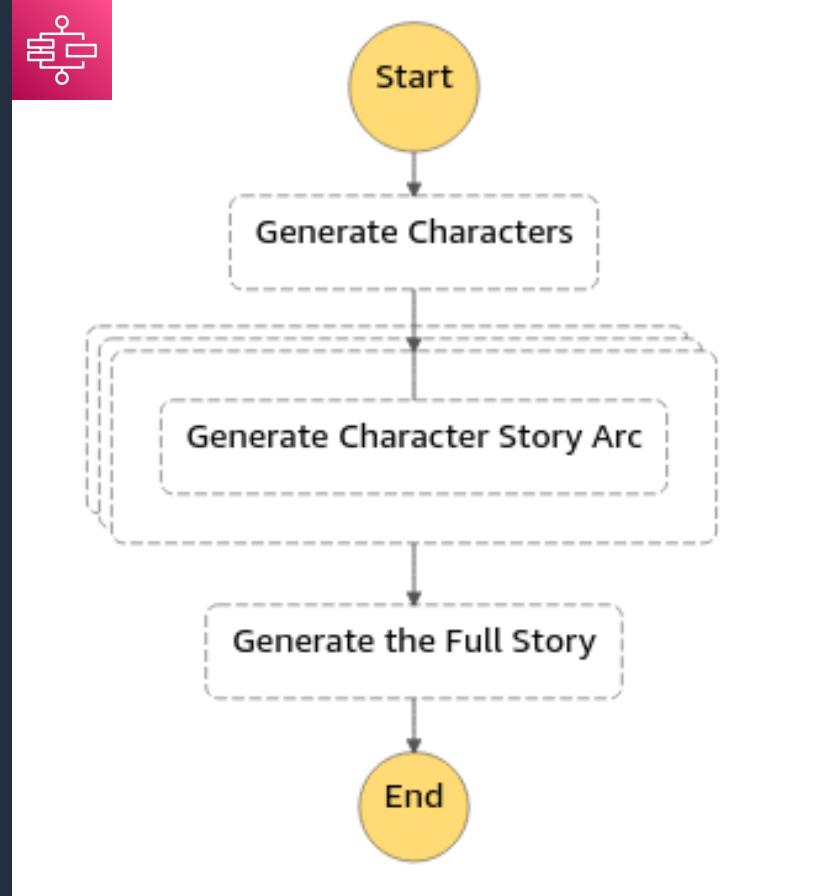
Use cases

- Domain specific Q/A
- Intelligent Document Processing
- Text Summarization

Transform data to vector embeddings



Orchestration of prompts



You are an award-winning fiction writer and you are writing a new story about {story_description}.

Before writing the story, describe five characters that will be in the story. Your response should be formatted as a JSON array, with each element in the array containing a "name" key for the character's name and a "description" key with the character's description.

An example of a valid response is below, inside ...



Amazon Bedrock
ANTHROP\C
Claude

FM Endpoints

- Amazon Bedrock
- Amazon SageMaker

Serverless Services

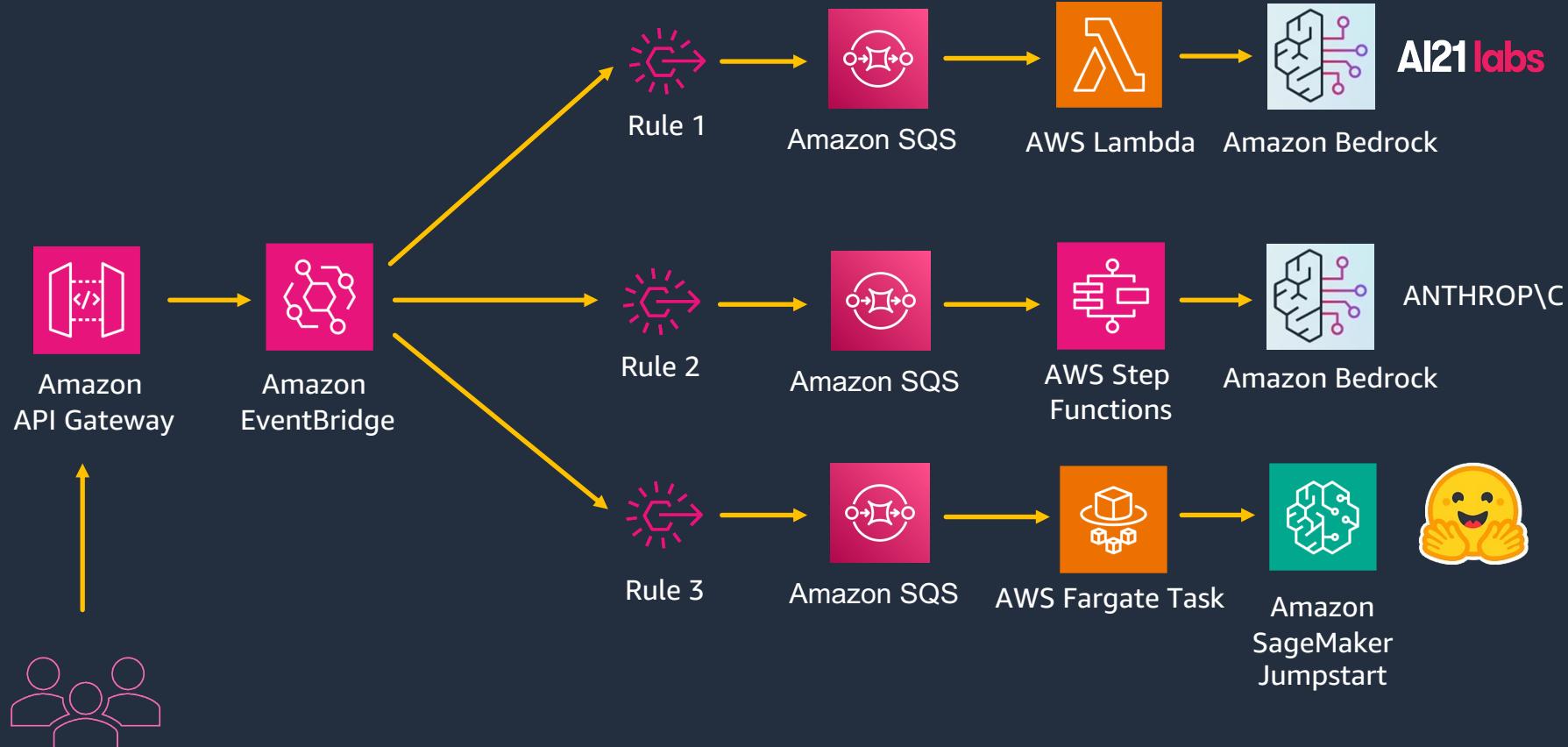
- Connect multiple prompts to generate complex content using Step Function
- Feed response from one model to the next

Use cases

- Writing blogs and articles
- Response validation
- Travel planning

Selecting model bases on **Type of Tasks**

Multiple models invocations by rules



FM Endpoint

- Amazon Bedrock
- Amazon SageMaker Jumpstart
- Models hosted on:
 - EKS
 - Other computes

Serverless Services

- Auth, Rate limiting, Caching with Amazon API Gateway
- Amazon EventBridge or Amazon SNS to fan-out
- Amazon SQS as subscribers
- AWS Lambda or AWS Step Functions getting from the queue

Use cases

- Multi-model evaluation
- Content generation from different models
- Selecting model for specific task

Simplify Workflow Automation using **Agents**

Workflow automation challenges



Knowledge workers
stretched, need
productivity tools



LLM's are powerful,
but they can't
take actions



Integration of databases
and systems is
expensive and slow



Building production
agents involves
complex engineering



Need diverse set of
programming languages
and interfaces

Agent basics

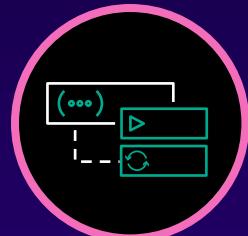


do this
for me...

Done. Here's
the result...

Agent

Instructions: "you are an agent that ..."



Actions

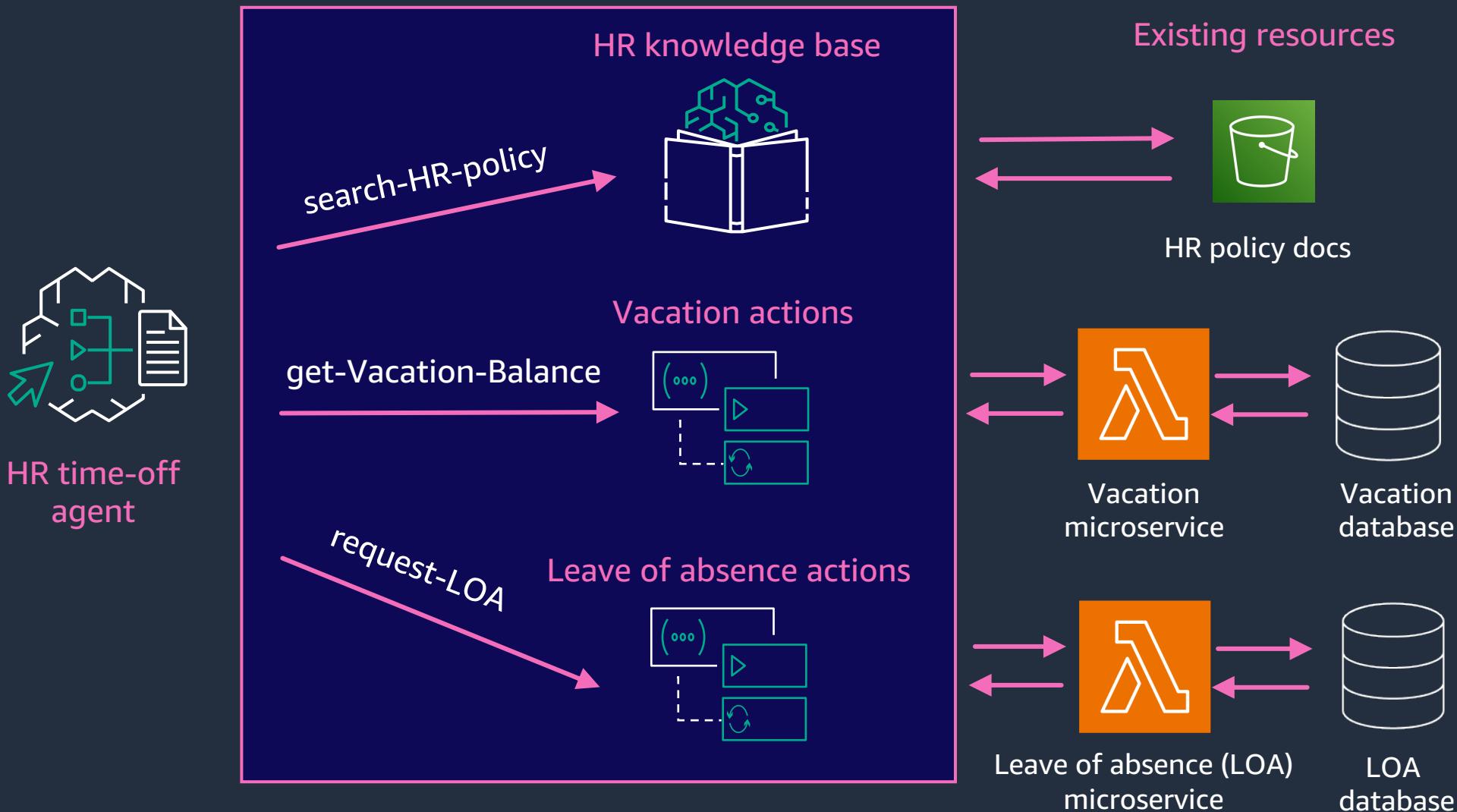


Knowledge Bases

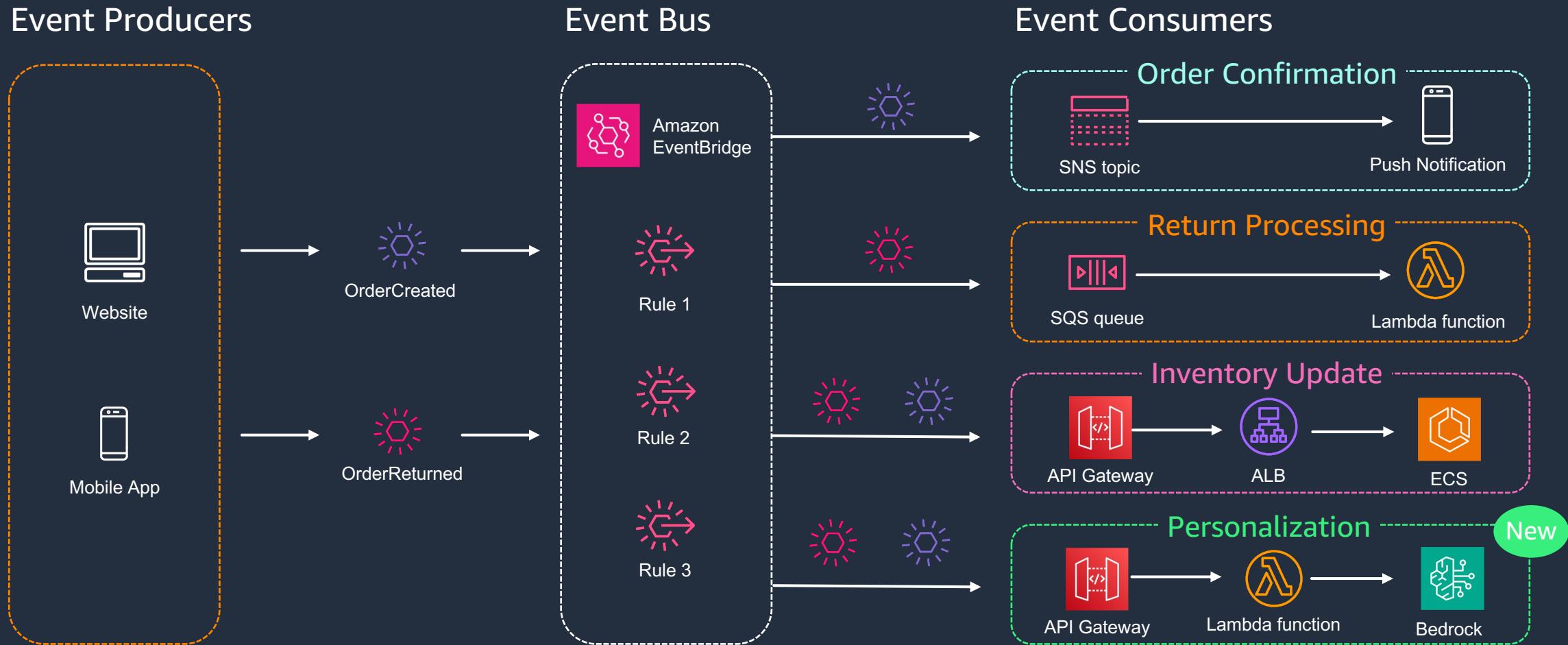


Amazon
Bedrock

Agents build on existing enterprise resources



Adding GenAI capabilities to e-commerce platform



Key takeaways

- Think asynchronously
- Think in events
- Think scale
- Build composable architectures



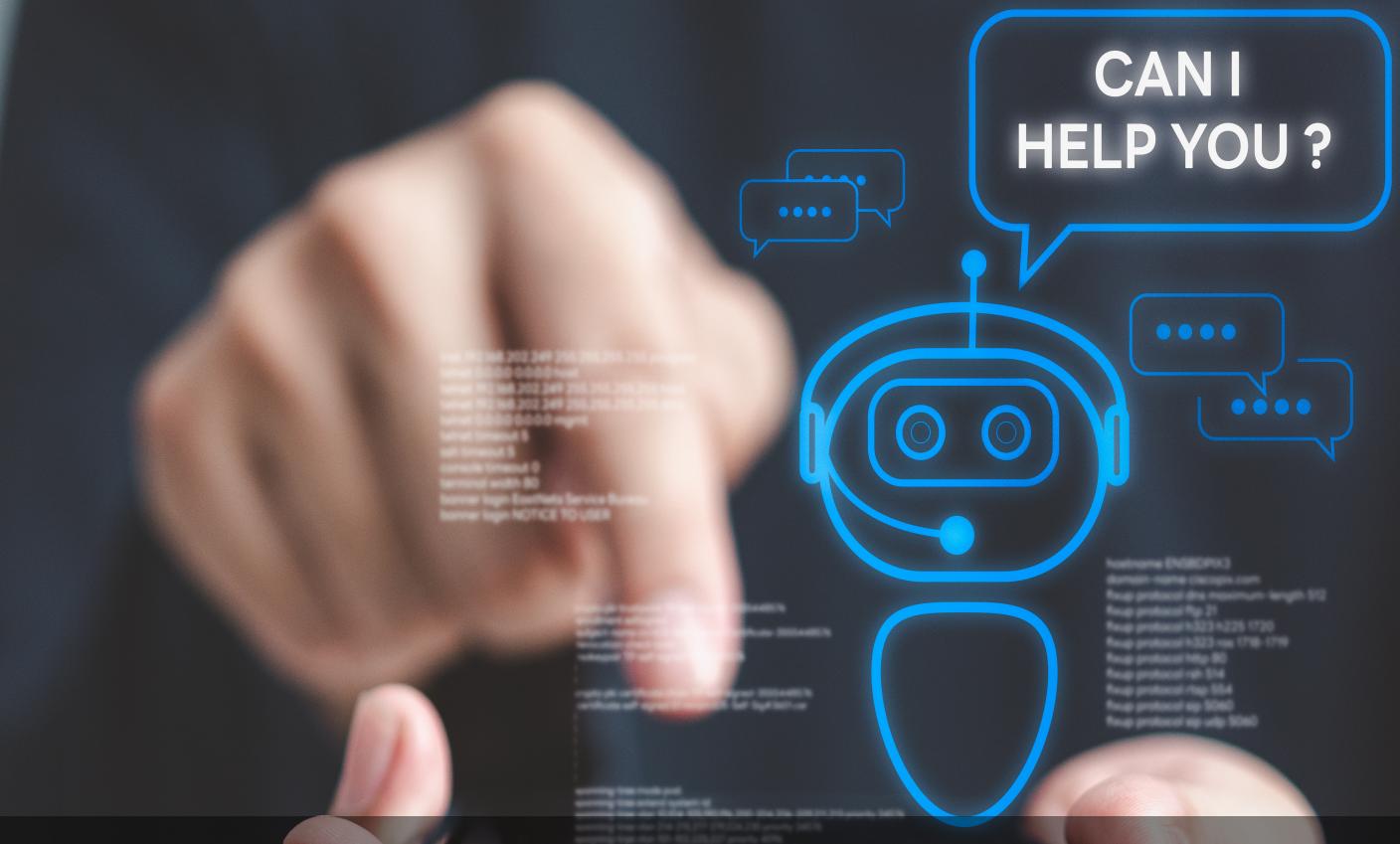


Thank you!

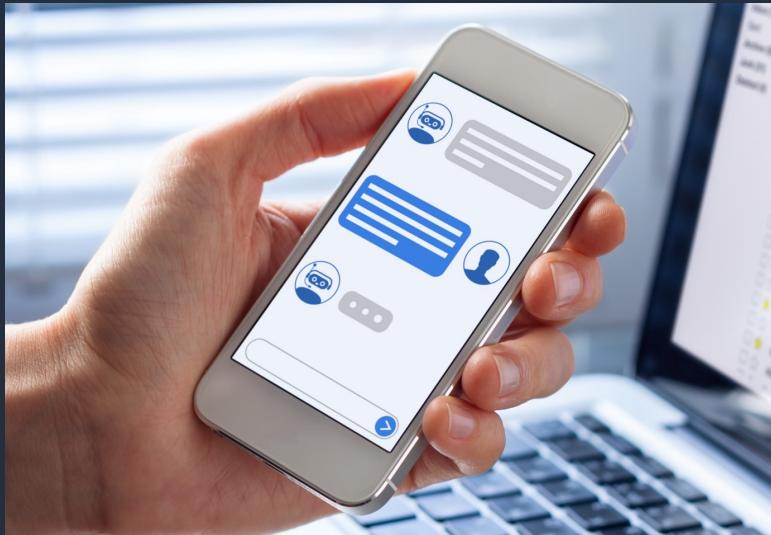


Use Cases

Use Case 1. Virtual Agent



Key capabilities



Self-service virtual agent

Frequently asked questions (Q&A)

Allow users to access FAQs quickly with voice/chat bots

Knowledge search

Allow users to find answers stored in disparate data sources



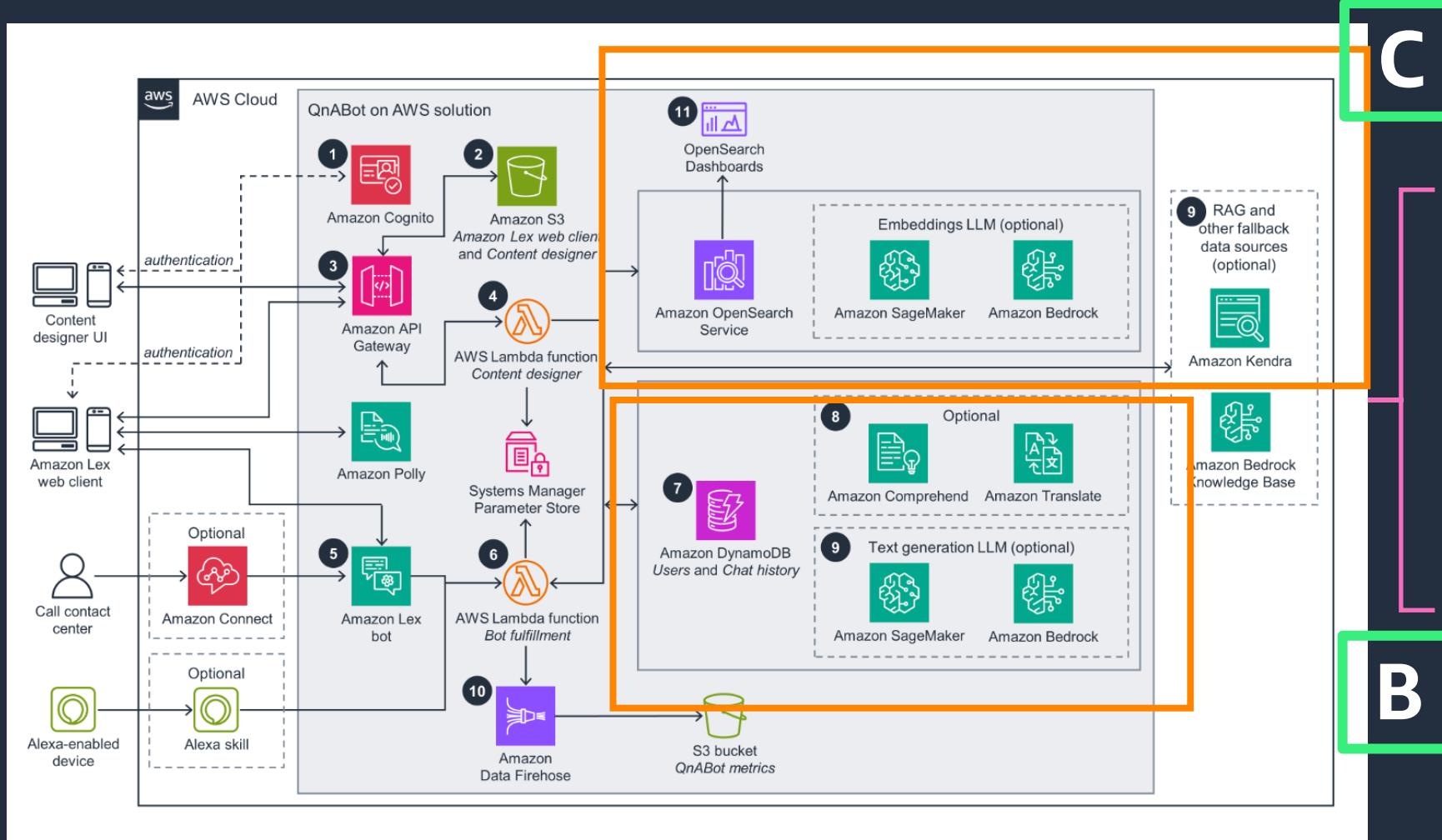
Transactions and task fulfillment

Automate simple tasks (i.e. bill payment, check balance, recommendation) with a sophisticated bot built with high quality ASR, NLU

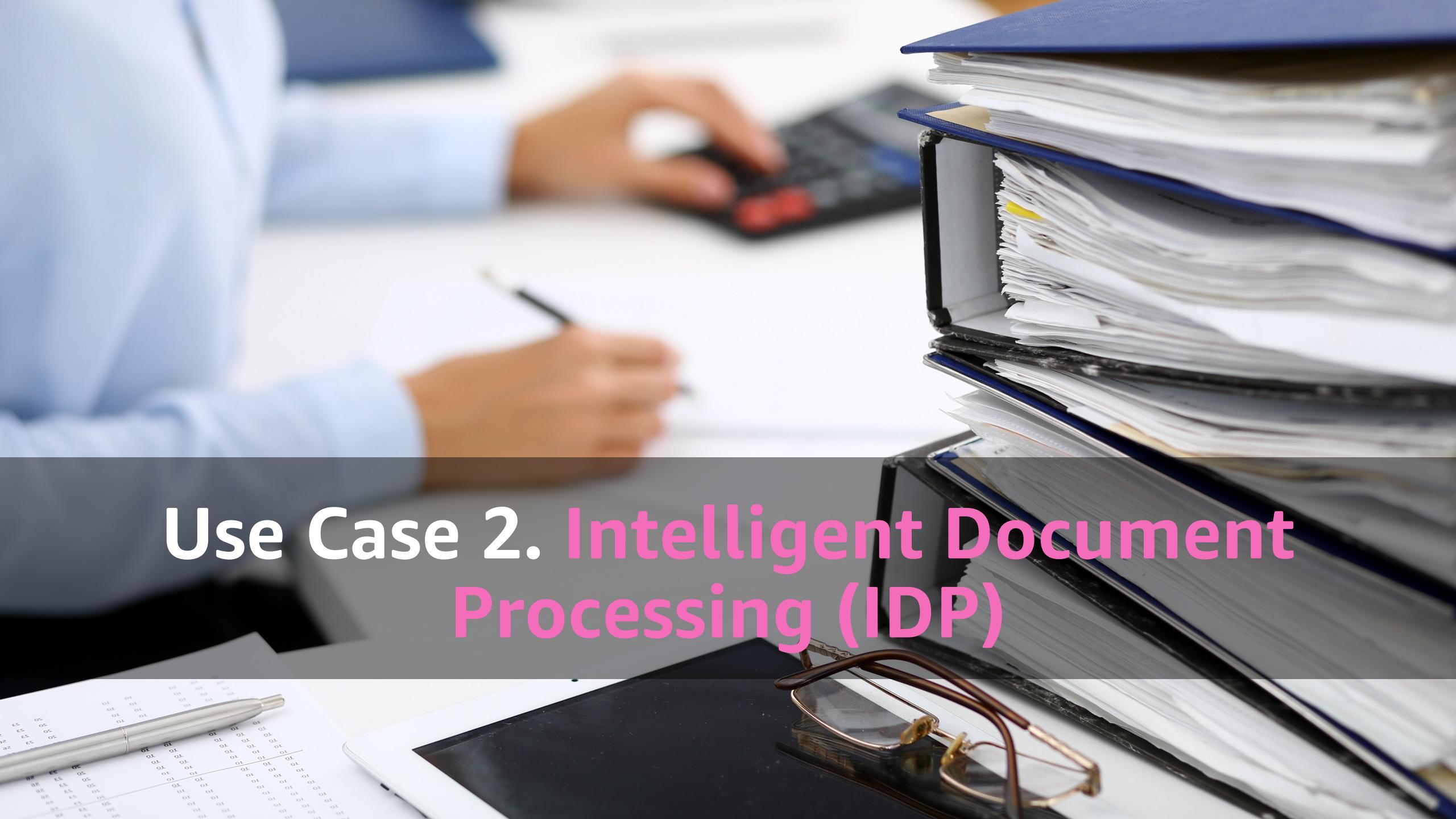
Intelligent routing

Capture information to assess which agent can best solve the user's request

Self service virtual agents solution overview



- Q & A chatbot
- Knowledge search
- Intelligent routing
- Fulfillments (i.e. bill payment, check balance, Next best actions)

A photograph of a person's hands working at a desk. One hand holds a calculator, and the other holds a pen over a piece of paper. In the foreground, there is a stack of papers, a pair of glasses, a tablet, and a pen. The background is blurred.

Use Case 2. Intelligent Document Processing (IDP)

Legacy document processes do not meet today's needs



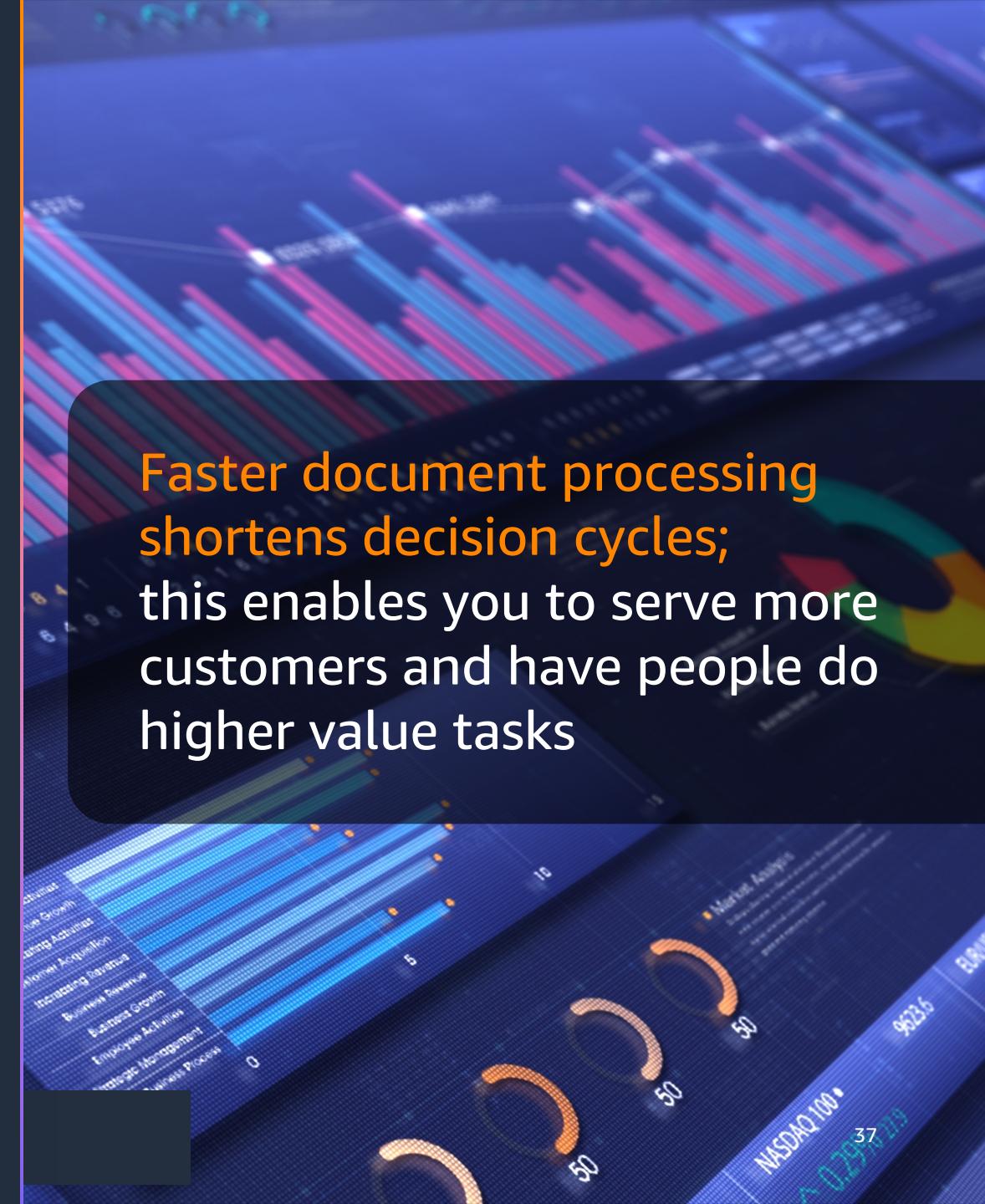
Legacy OCR and manual processes are time-consuming, error-prone, and expensive



Manual processes do not scale easily with document volume



Difficult to find useful information needed for business decisions



Faster document processing shortens decision cycles; this enables you to serve more customers and have people do higher value tasks

IDP at AWS

Ready to use AI tools that help you classify and extract key data from documents

Ready to use AIs that classify and extract key business data from your documents



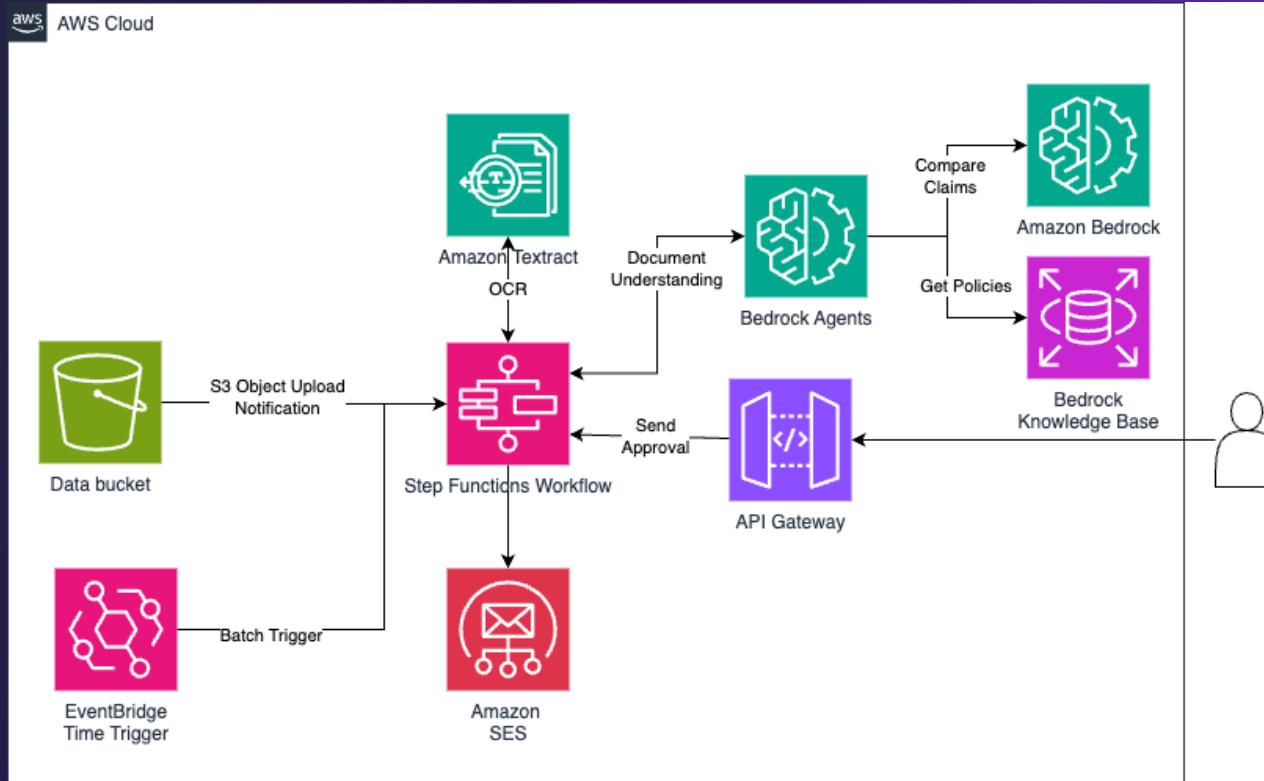
Enrich and validate key insights before sending to downstream systems

Serve end customers faster

Reduce the total cost of document processing

Intelligent Document Processing (IDP) workflow overview

Claims processing with multilingual support

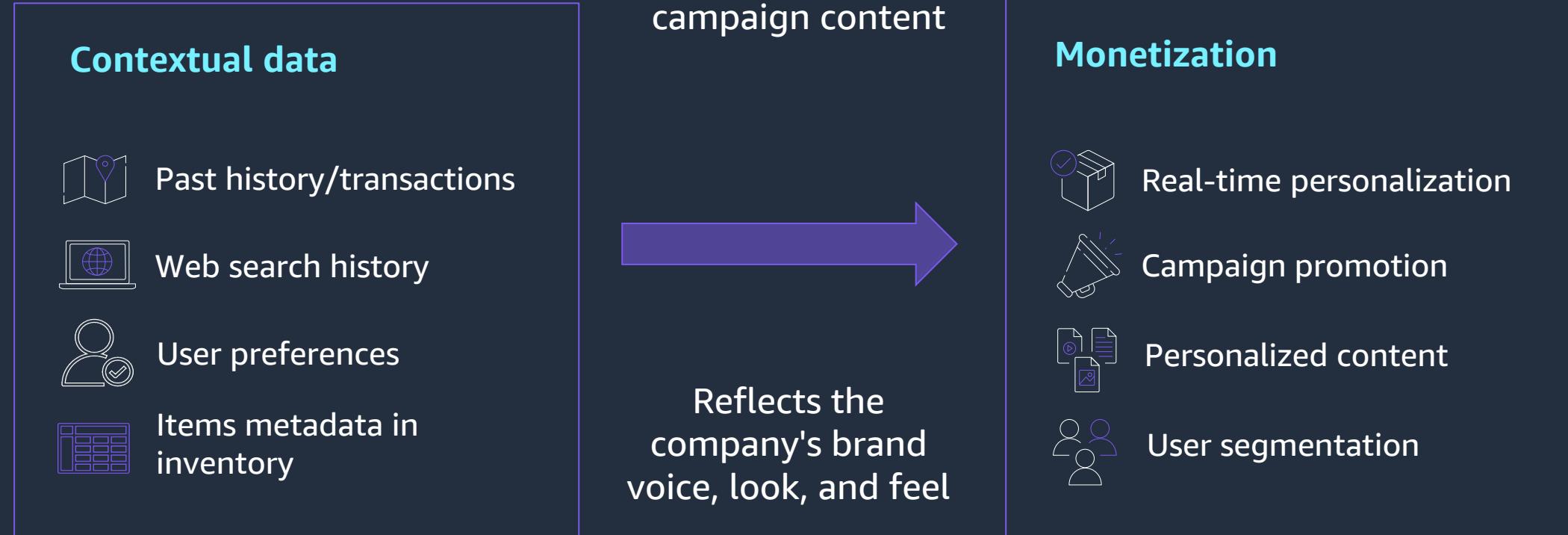


- Preprocessing with OCR
- Knowledge search
- Agents based Intelligent routing
- Fulfillments (i.e. notification, approval)

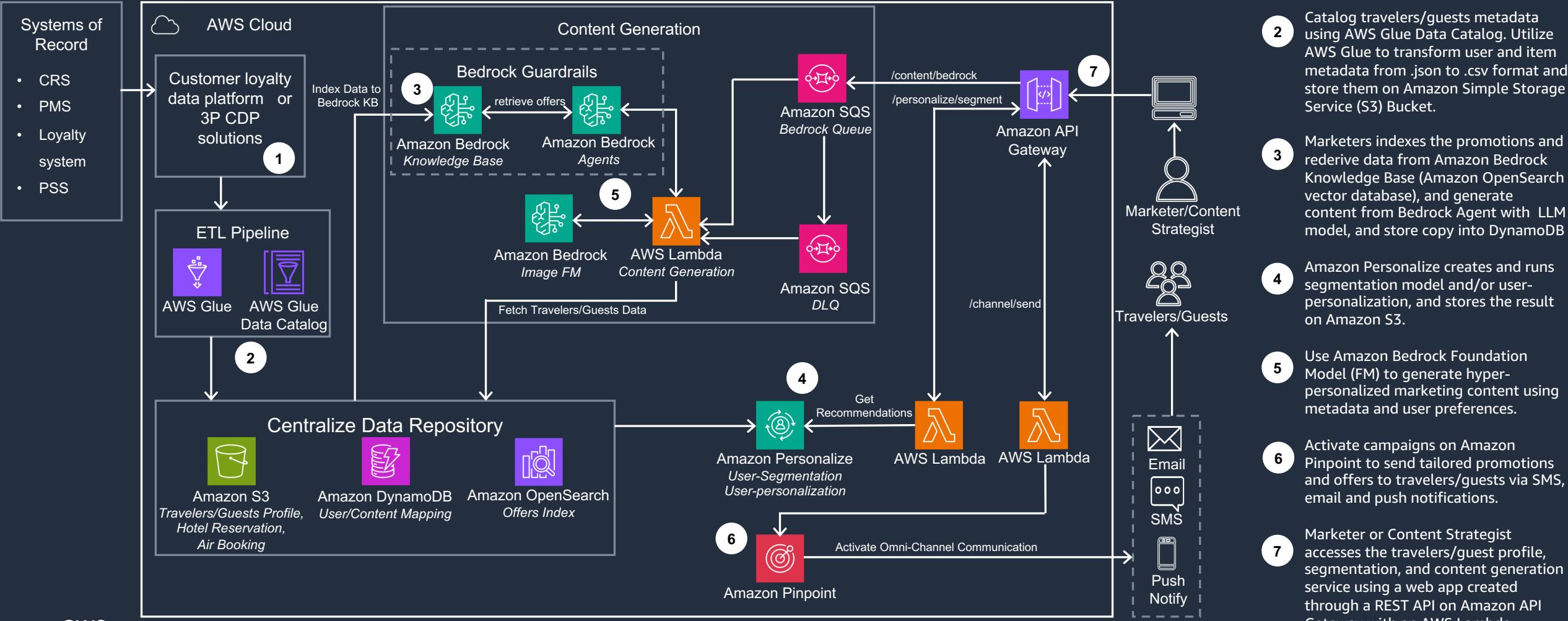
Use Case 3. Hyper-personalization



Key capabilities



Hyper-personalization Reference Architecture



Why Serverless

New AppDev norm for innovators & disrupted industries



Faster Time to Business Value
Daily to multiple times per day vs. Monthly on VMs



Cost Savings
Non-expiry Free Tier ; Allocates resources according to demand



Culture of Innovation
40% productivity increase and more focus on innovation



Excellent Scalability, Security, and Availability
Build-in reliability and security; automatic scaling.

Proven Path



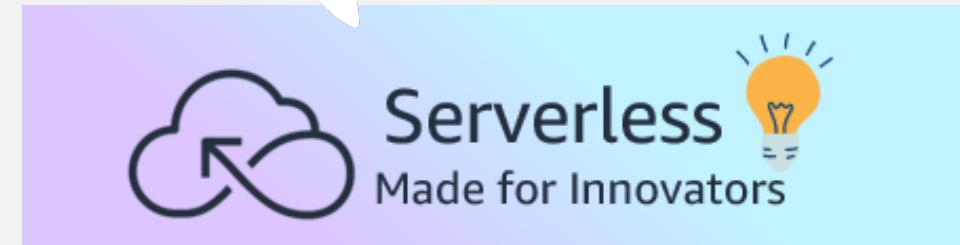
Serverless **simplifies** modernization



... and helps you fast track to **Innovate** phase

AWS Support

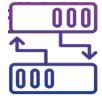
1. Partner Network & Funding
2. Experience-based Acceleration Workshop
3. Professional Services
4. Training & Certification
5. Hands-on immersion day



Use Case



IT Automation



Data Processing



Event-driven Architecture



Web Application



AI / Machine Learning



Mobile Application

1,700+ Serverless customers run on AWS

1. Cost saving up to 57% when use in right case
2. 30-40% productivity increase vs. on-prem
3. Faster time-to-market (daily to multiple times per day vs. fortnightly on VMs)

A Perfect Choice made for Innovators

Customers



BOSCH



THE SHONET



SIEMENS

AUTODESK



7-ELEVEN

Katalon



Expedia®

coinbase

AlteaCare

JUNIPER NETWORKS



Foundation Model Builders & Providers

As a Service



Bedrock
(announced)



Proprietary

co:here

Light^{on}

AI21^{labs}

Public

stability.ai

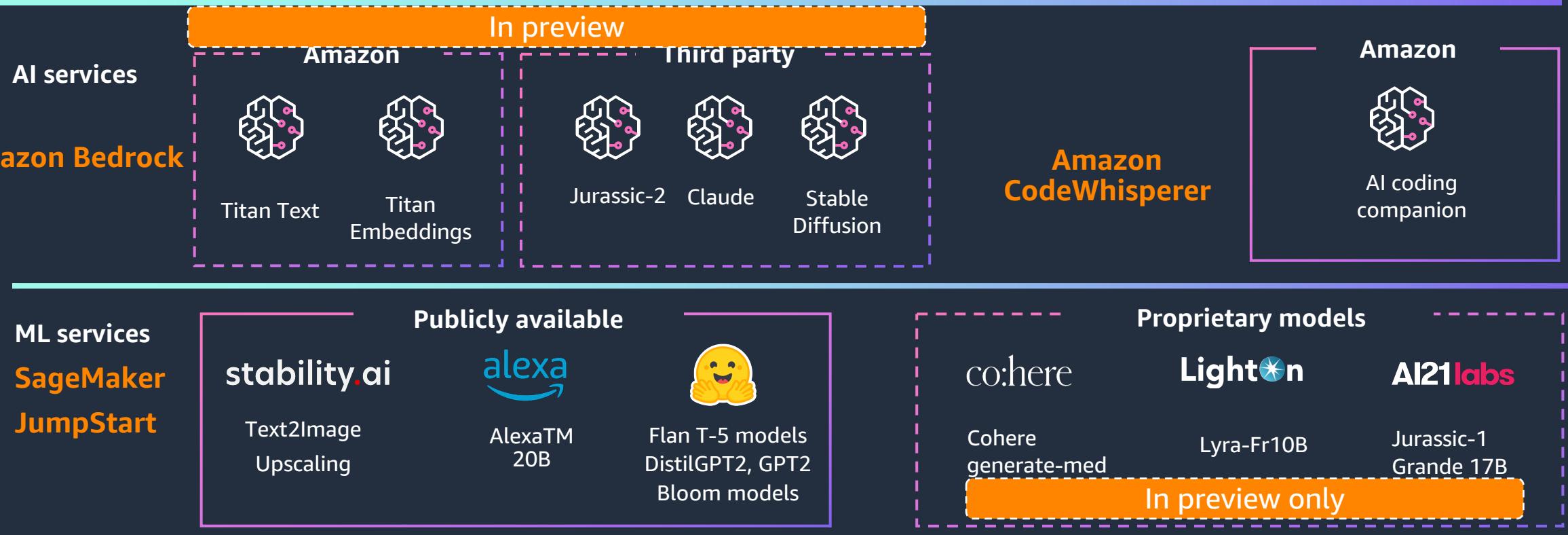
alexa^{TM 20B}



HuggingFace

Amazon Generative AI portfolio

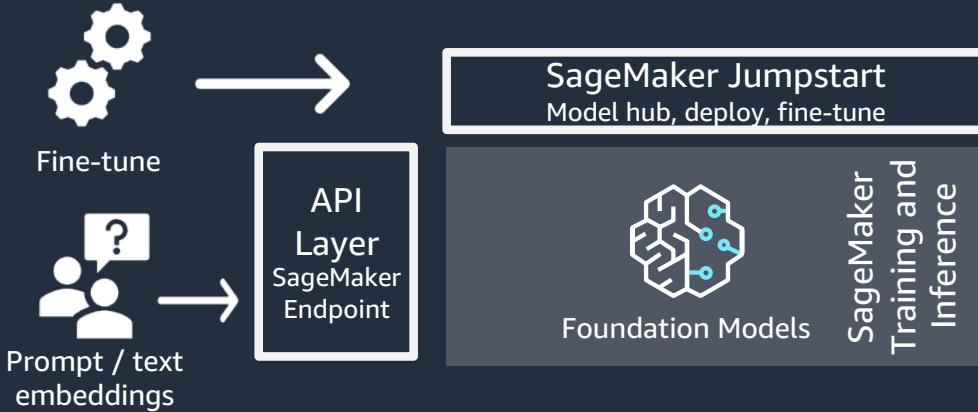
CHOICE OF MANY FOUNDATION MODELS



© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved. Amazon Confidential and Trademark.

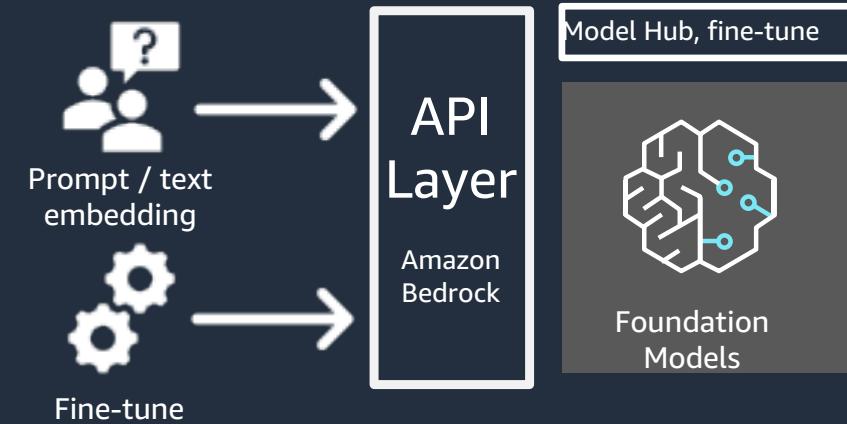
How do I access foundation models?

Amazon SageMaker JumpStart



- Select the model: ML hub with FMs, built-in algorithms, and prebuilt ML solutions that you can deploy with just a few clicks
- Configuration based deployment: Deploy FM as SageMaker endpoint (hosting)
- Invoke Endpoint
- Managed infrastructure

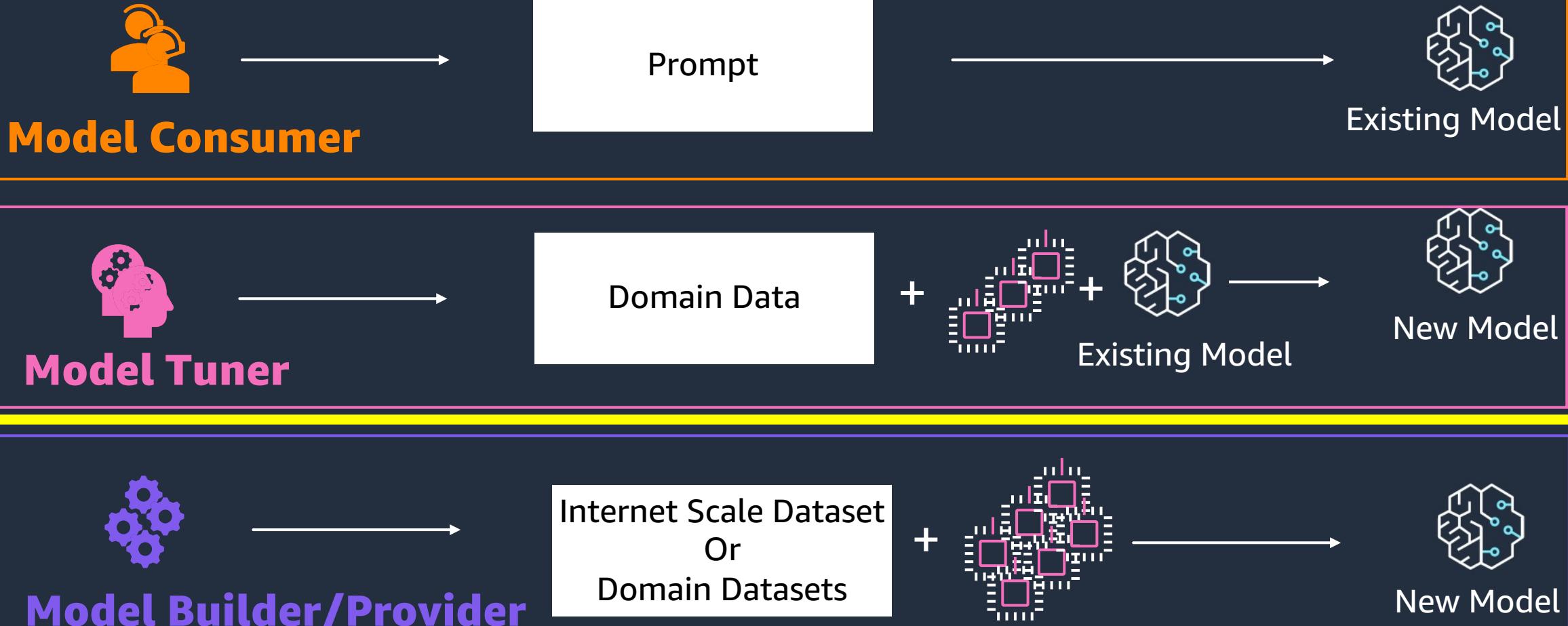
Amazon Bedrock



- Select the model: The easiest way to build and scale generative AI applications with FMs
- -
- Invoke Endpoint: Access directly or fine-tune foundation model using API
- Serverless

Gen AI Integration with Serverless

Serverless



Build value at scale and speed with Gen AI



Generative AI Application:
Business Value

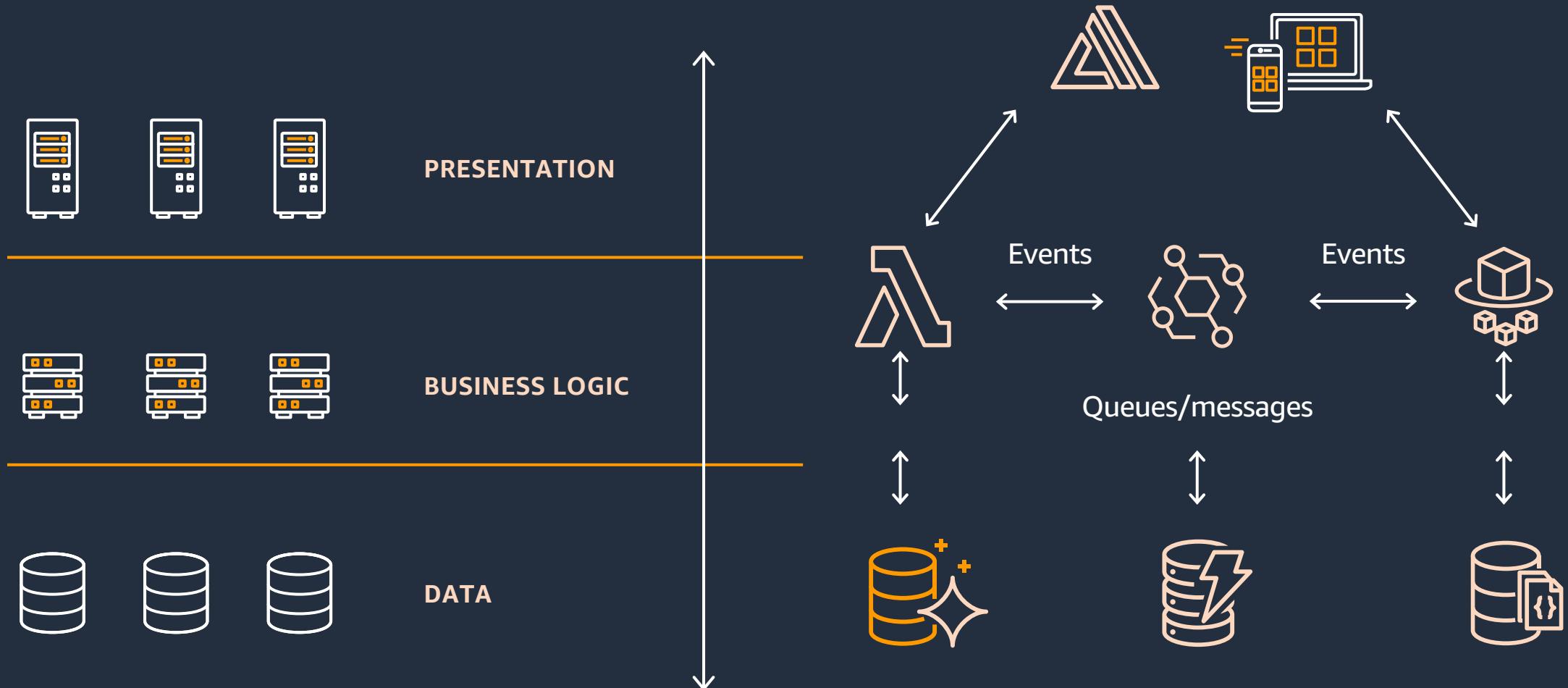
Modern Application:
**BUILT ON
SERVERLESS**

PROCESS: GOVERNANCE
& SECURITY

MINDSET: EXPERIMENTAL &
CUSTOMER OBSESSED

PEOPLE: NEW SKILLS &
ROLES

Key characteristic of Serverless: event-driven microservices



Combining Speed and Power

Speed of Serverless



AWS
Lambda



Amazon
ECS



AWS
Fargate



AWS
Step
Functions



Amazon
EventBridge

Power of Gen AI



Amazon
SageMaker

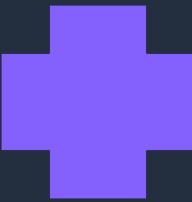


Amazon
Bedrock



Amazon
Code Whisperer

Rapid delivery of
smarter
applications and
features with
focus on
Innovation



Generate personalized electronic direct mail (EDM) for the new car model promotion

Input data

 Contextual metadata

 Augmented metadata

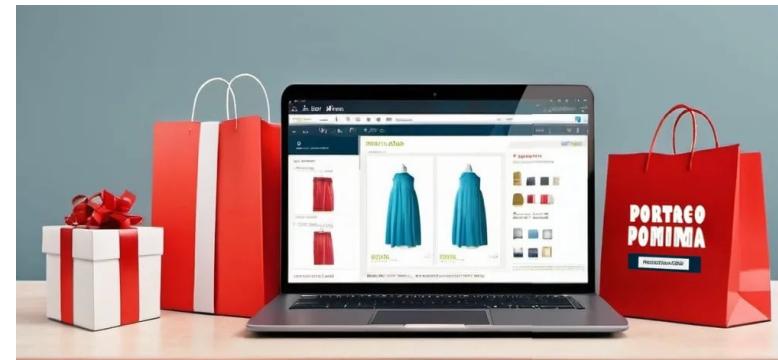
 Campaign questions

Output contents

 Email title/
body/banner

 Email banner
image

 Fulfilment/next best
action



Subject: Upgrade Your Tech with Our Top-Rated Laptops

Dear [Customer Name],

As a valued customer, we're excited to share with you our latest collection of high-performance laptops, tailored to meet your unique needs and preferences. At [Company Name], we understand the importance of technology in your daily life, and we're committed to providing you with the best products and experiences.

Introducing the TechBrand Laptop: A Powerhouse of Performance

Our featured product, the TechBrand Laptop, is a true masterpiece of engineering. With its sleek silver aluminum body and 15-inch touchscreen display, this laptop is not only visually stunning but also packed with cutting-edge features. Boasting a powerful Intel i7 processor, 16GB of RAM, and a lightning-fast 512GB SSD, the TechBrand Laptop is designed to handle even the most demanding tasks with ease.

Personalized for Your Lifestyle

Based on your interests in technology and active lifestyle, we believe the TechBrand Laptop is the perfect companion for your on-the-go adventures. Whether you're working remotely, streaming your favorite shows, or editing high-resolution videos, this laptop delivers unparalleled performance and portability.

Exclusive Offer for [Loyalty Status] Members

As a valued [Loyalty Status] member, we're thrilled to offer you an exclusive discount on the TechBrand Laptop. For a limited time, you can enjoy [Discount Percentage] off the regular price, making this top-notch device even more accessible.

Don't miss out on this incredible opportunity to elevate your tech game. Visit our website or your nearest [Company Name] store today to experience the TechBrand Laptop firsthand and take advantage of this exclusive offer.

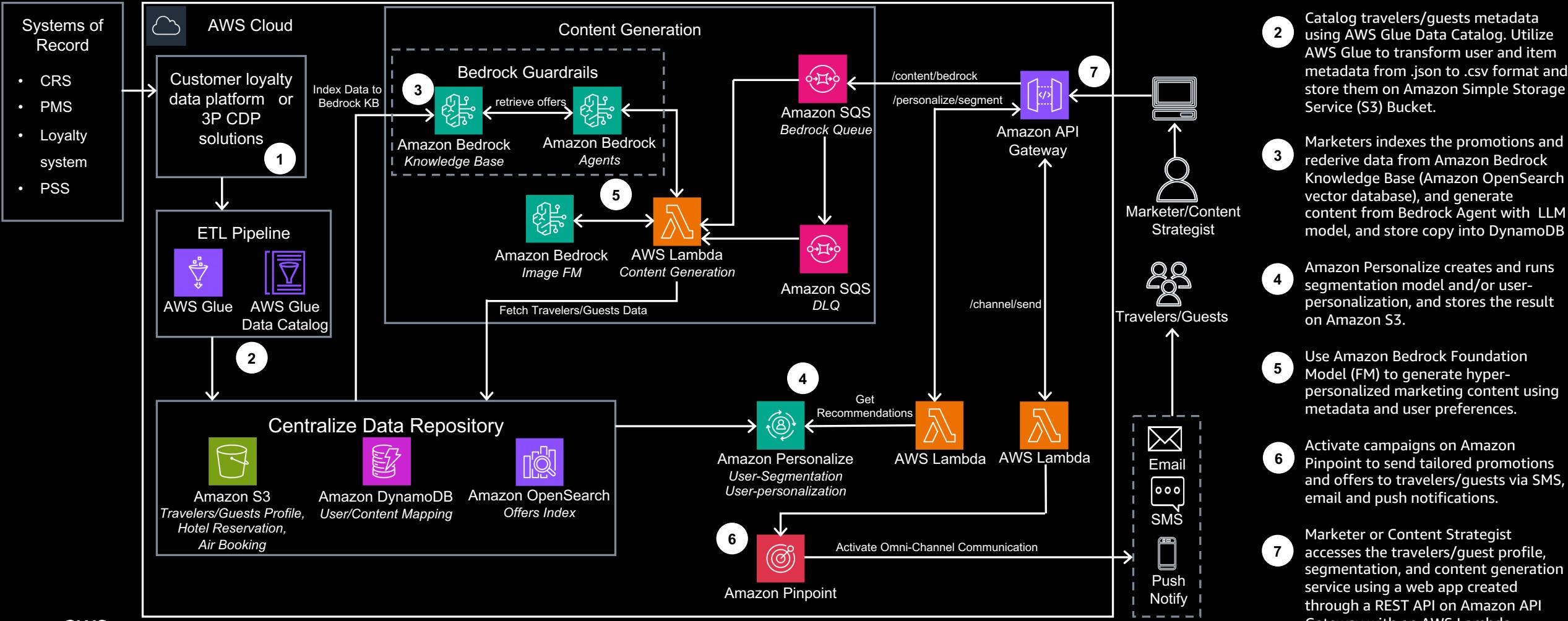
Thank you for your continued loyalty and trust in [Company Name]. We look forward to providing you with an exceptional shopping experience.

Best regards,
[Your Name]
Chief Marketing Officer
[Company Name]

Generated by LLM



Hyper-personalization Reference Architecture



Demo – Digital Banking hyper-personalization

The screenshot shows a web browser window for 'Agenticflow AI' with the URL app.agenticflow.ai/dashboard. The left sidebar contains navigation links for Teamspaces (FSI), AI Tools (Chat, Assistant), Workforce (Workflows, Templates, Assets, History), Monitor, and Activity center. A user profile for 'joseph_tran' is at the bottom. The main content area displays a bar chart titled 'Total INTEREST RATE by PRODUCT NAME'. The chart lists various banking products and their interest rates:

PRODUCT NAME	INTEREST RATE
Secured Credit Card	24.99
Travel Rewards C...	18.99
Cash Back Reward...	16.99
Debt Consolidati...	9.99
Small Business L...	8.5
Home Improveme...	7.99
Home Equity Line...	5.5
Federal Student ...	4.99
30-Year Fixed Ra...	4.25
New Car Loan	3.99

A scrollable list of additional product names is shown on the right side of the chart area.