ZURICH | September 4, 2024

aws SUMMIT

AIM303

# Deep Dive: Building AI agents using Amazon Bedrock

**Viktor Vedmich**
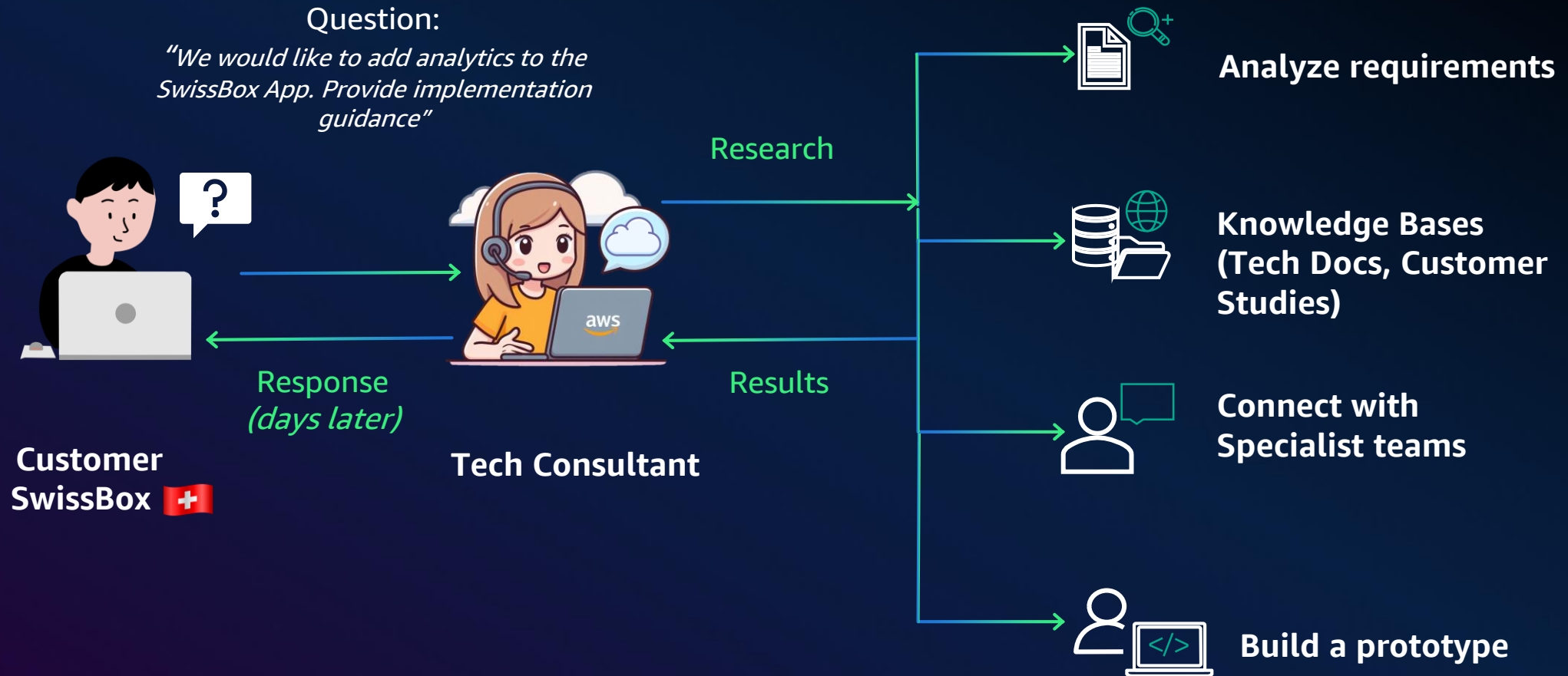
Senior Developer Advocate
AWS

**Viktoria Semaan**

Senior Developer Advocate
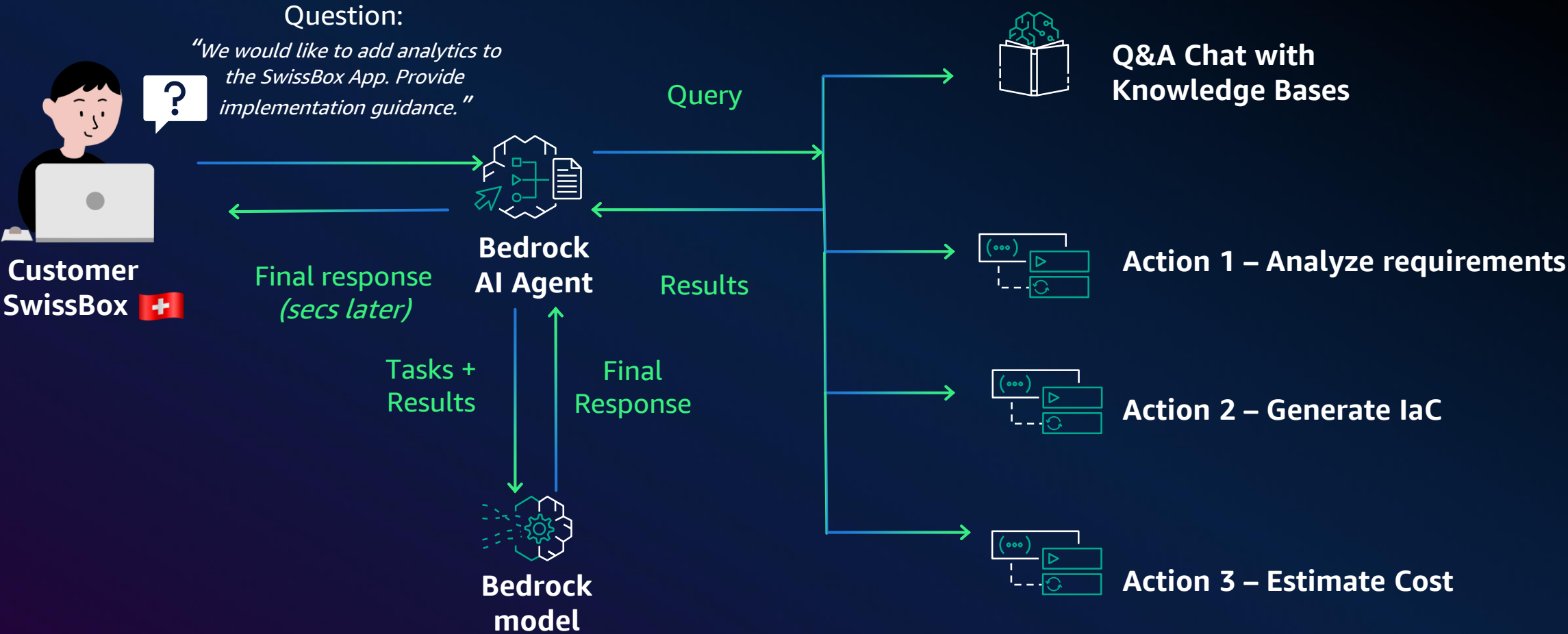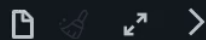AWS

# Human Workflow

# AI Agent Workflow

Question:
*"We would like to add analytics to the SwissBox App. Provide implementation guidance."*

**Customer SwissBox** 🇨🇭

Query

**Bedrock AI Agent**

Final response *(secs later)*

Results

Tasks + Results

Final Response

**Bedrock model**

**Q&A Chat with Knowledge Bases**

**Action 1 – Analyze requirements**

**Action 2 – Generate IaC**

**Action 3 – Estimate Cost**

strategy for the Redshift tables based on the data characteristics and query patterns. - Optionally, use Redshift's massively parallel processing (MPP) architecture to run complex analytical queries on the data. 4. Analytics and Visualization with QuickSight: - Connect Amazon QuickSight to the Redshift cluster as a data source. - Use QuickSight to create visualizations, dashboards, and reports on the data stored in Redshift. - Share and publish the QuickSight dashboards with SwissBox users. - Optionally, use QuickSight's machine learning capabilities for forecasting and anomaly detection. To implement this solution, I can provide sample Terraform code for provisioning the required AWS services and resources. Let me know if you need any clarification or have additional requirements.

Show trace ›

aws_iam_role.glue_role.arn job_temp_bucket = aws_s3_bucket.glue_temp_bucket.id } # Redshift module module "redshift" { source = "./modules/redshift" cluster_identifier = "swissbox-redshift-cluster" database_name = "swissboxdb" node_type = "dc2.large" number_of_nodes = 2 username = "swissboxadmin" # Other Redshift config } # QuickSight module module "quicksight" { source = "./modules/quicksight" quicksight_user_name = "swissbox_analyst" redshift_data_source_name = "swissbox-redshift-data" redshift_cluster_id = module.redshift.redshift_cluster_id redshift_database = "swissboxdb" redshift_user = "swissboxadmin" # Other QuickSight config } # Detailed module code would go into respective module directories

Show trace ›

Enter your message here

Run

AMAZON BEDROCK

Models to choose
Model Evaluation

INTRO
Use Case

# Amazon Bedrock Foundation Models

| AI21 labs | amazon | ANTHROP\C | cohere | Meta | MISTRAL AI_ | stability.ai |
|---|---|---|---|---|---|---|
| Contextual answers, summarization, paraphrasing | Text summarization, generation, Q&A, search, image generation | Summarization, complex reasoning, writing, coding | Text generation, search, classification | Q&A and reading comprehension | Text summarization, text classification, text completion, code generation, Q&A | High-quality images and art |
| Jamba-Instruct | Amazon Titan Text Premier | Claude 3.5 Sonnet | Command | Llama 3 8B | Mistral Small | Stable Diffusion XL1.0 |
| Jurassic-2 Ultra | Amazon Titan Text Lite | Claude 3 Opus | Command Light | Llama 3 70B | Mistral Large | Stable Diffusion XL 0.8 |
| Jurassic-2 Mid | Amazon Titan Text Express | Claude 3 Sonnet | Embed English | Llama 3.1 8B | Mistral 7B | |
| | Amazon Titan Text Embeddings | Claude 3 Haiku | Embed Multilingual | Llama 3.1 70B | Mixtral 8x7B | |
| | Amazon Titan Text Embeddings V2 | Claude 2.1 | Command R+ | Llama 3.1 405B | | |
| | Amazon Titan Multimodal Embeddings | Claude 2 | Command R | | | |
| | Amazon Titan Image Generator | Claude Instant | | | | |

AMAZON BEDROCK

Models to choose
Model Evaluation

CUSTOMIZATION

Prompt Engineering
Knowledge Base (RAG)

INTRO
Use Case

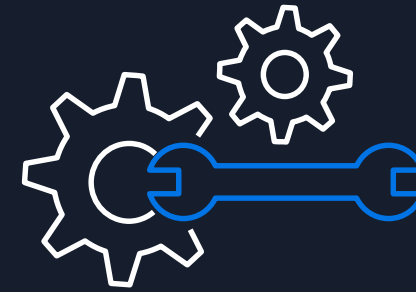# Customizing model responses for your business

## Fine-tuning

**PURPOSE**

Maximizing accuracy for **specific tasks**

**DATA NEED**

**Small number** of labeled examples

## Continued pretraining

**PURPOSE**

Maintaining model accuracy for **your domain**

**DATA NEED**

**Large number** of unlabeled datasets

# Knowledge Base: End-to-End RAG



**Text Generation Workflow**

User
User Input
Embeddings model
Embedding
0.89 | -0.02 | -0.53 | 0.95 | 0.17 | ••• | -0.38
Context
Prompt augmentation
Large Language Model
Response

**Data Ingestion Workflow**

Semantic search
Vector store
Embeddings model
Document chunks
Data source

AGENT TOOLS

SA Q&A
Generate IaC
Cost Estimation

AMAZON BEDROCK

Models to choose
Model Evaluation

CUSTOMIZATION

Prompt Engineering
Knowledge Base (RAG)

INTRO
Use Case

# AI Agent: Detailed Flow

**Task:**

*How to build analytics for the SwissBox app. How much it will cost to run it monthly?*

**Bedrock Agent**

**PROMPT**

Conversation history

Actions, Knowledge Bases

Instructions

Task

**Bedrock Model**

Chain of thought

Step 1

Step 2

…

Step n

Search

Results

**Knowledge Bases**

API call

Results

**Action Groups**

Results

**Final response**

Task + results

**Bedrock model**

Final response

*Configuration steps: 1/…..*

*Expected cost: $178*

Decompose into steps using KBs and actions

↓

Execute action or search knowledge base

↓

Observe results

↓

Think about next step

# Tool 1: Question & Answer Bot



Solutions Architect Agent

Tool 1: SA Q&A

How to…?

Tool 2: Generate IaC

Tool 3: Estimate Cost

AWS documentation

# Knowledge Bases

**Data sources**

Amazon S3

Web Crawler

Atlassian Confluence

Salesforce

Microsoft SharePoint

**Embeddings model**

Titan Text Embeddings v2
By Amazon

Embed English v3
By Cohere
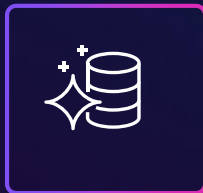
Titan Text Embeddings G1 – Text v1.2
By Amazon

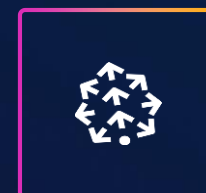Embed Multilingual v3
By Cohere

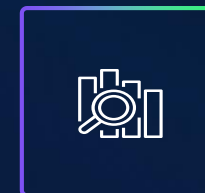**Vector database**
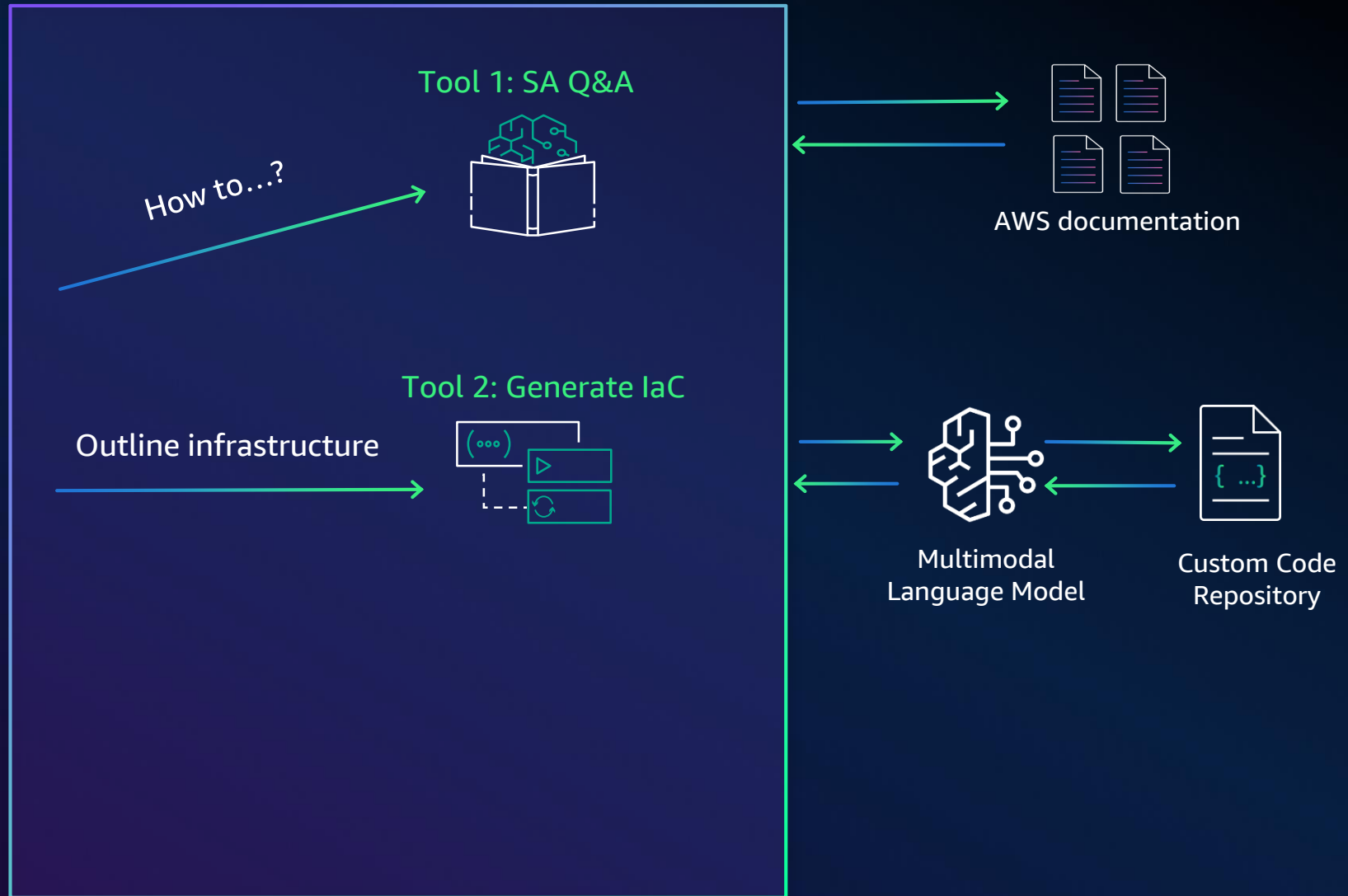
Amazon Aurora
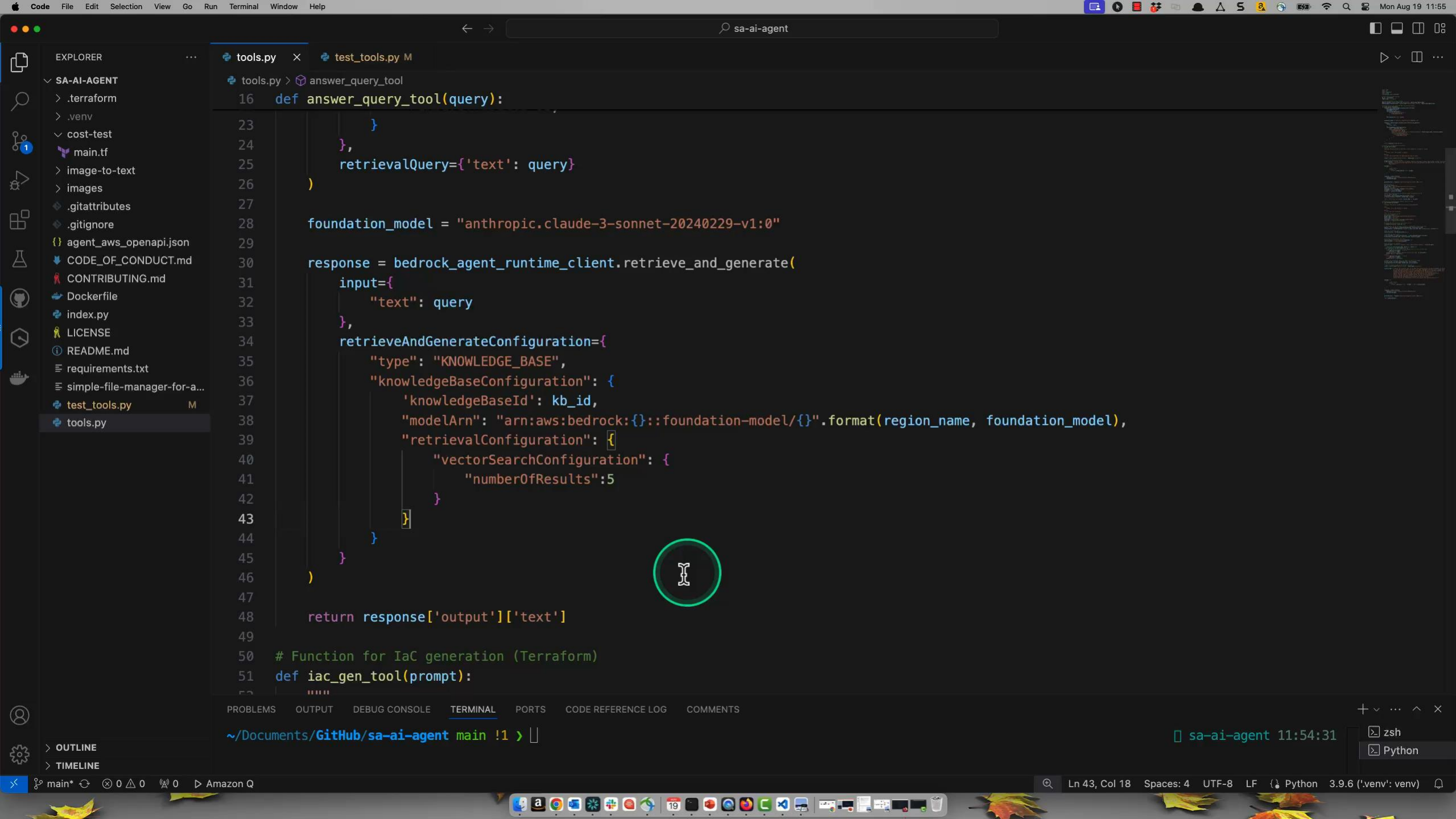
Redis Enterprise Cloud

MongoDB Atlas

Pinecone

Vector Engine For Amazon OpenSearch Serverless

# Tool 2: Generate IaC



Tool 1: SA Q&A

How to…?

AWS documentation

Solutions Architect Agent

Tool 2: Generate IaC

Outline infrastructure

Multimodal Language Model

Custom Code Repository

```python
def answer_query_tool(query):
            }
        },
        retrievalQuery={'text': query}
    )

    foundation_model = "anthropic.claude-3-sonnet-20240229-v1:0"

    response = bedrock_agent_runtime_client.retrieve_and_generate(
        input={
            "text": query
        },
        retrieveAndGenerateConfiguration={
            "type": "KNOWLEDGE_BASE",
            "knowledgeBaseConfiguration": {
                'knowledgeBaseId': kb_id,
                "modelArn": "arn:aws:bedrock:{}::foundation-model/{}".format(region_name, foundation_model),
                "retrievalConfiguration": {
                    "vectorSearchConfiguration": {
                        "numberOfResults":5
                    }
                }
            }
        }
    )

    return response['output']['text']

# Function for IaC generation (Terraform)
def iac_gen_tool(prompt):
    """
```

```
~/Documents/GitHub/sa-ai-agent main !1 >
```

# Tool 3: Estimate Cost

Tool 1: SA Q&A

How to…?

Solutions Architect Agent

Outline infrastructure

Tool 2: Generate IaC

Estimate cost

Tool 3: Estimate Cost
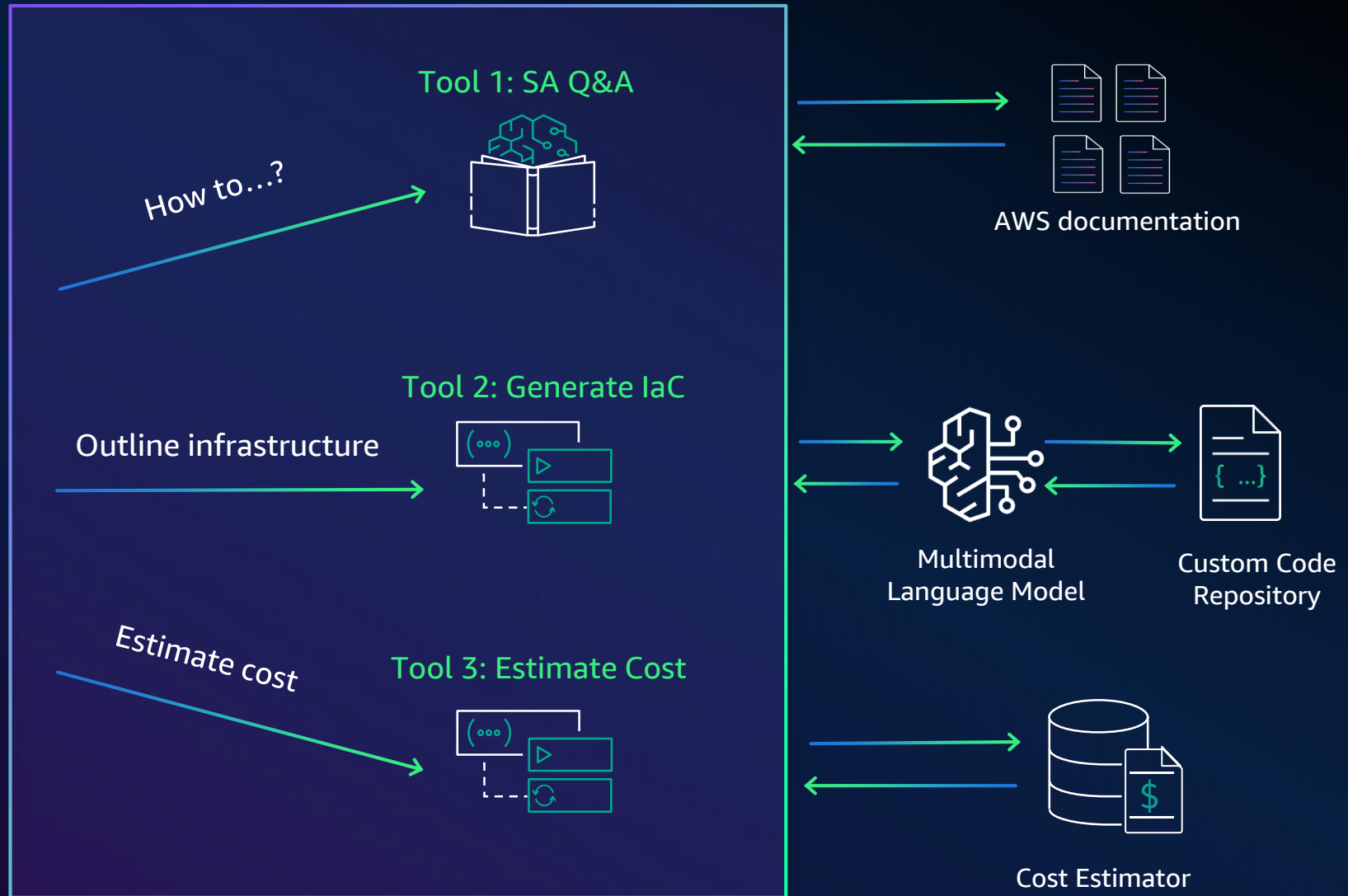
AWS documentation
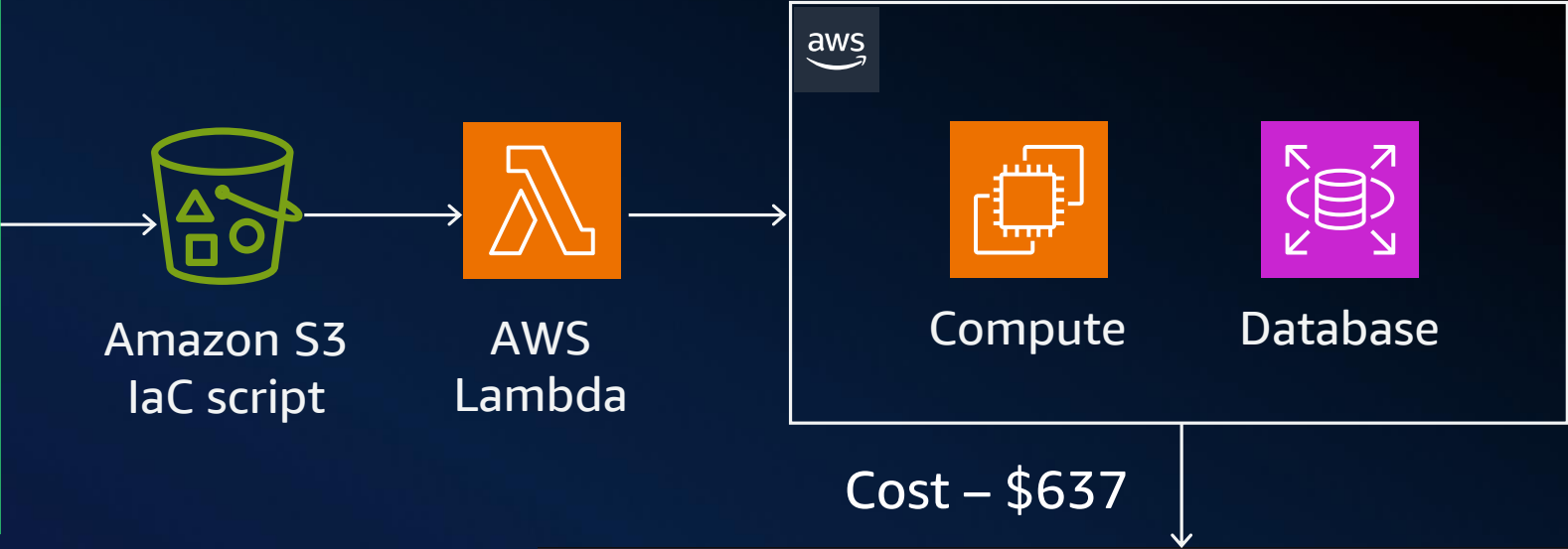
Multimodal Language Model

Custom Code Repository

Cost Estimator

# Cost Estimator: Configuration

```
resource "aws_db_instance" "db" {
  allocated_storage = 100
  instance_class    = db.r5.large
  ...
}

resource "aws_instance" "ec2" {
  count         = 2
  instance_type = m5.xlarge
  ...
}
```



Amazon S3
IaC script

AWS
Lambda

Compute          Database

Cost – $637

```
Project: project_t3kyto6a

Name                                                              Monthly Qty  Unit        Monthly Cost

aws_db_instance.rds_instance
├─ Database instance (on-demand, Single-AZ, db.r5.large)              730      hours            $175.20
├─ Storage (general purpose SSD, gp2)                                 100      GB                $11.50
├─ Additional backup storage                                         238      GB                $22.61  *
└─ Extended support (year 1)                                       1,460      vCPU-hours       $146.00

aws_instance.ec2_instances[0]
├─ Instance usage (Linux/UNIX, on-demand, m5.xlarge)                 730      hours            $140.16
└─ root_block_device
   └─ Storage (general purpose SSD, gp2)                               8      GB

aws_instance.ec2_instances[1]
├─ Instance usage (Linux/UNIX, on-demand, m5.xlarge)                 730      hou
└─ root_block_device
   └─ Storage (general purpose SSD, gp2)                               8      G

OVERALL TOTAL

*Usage costs were estimated with usage defaults from Infracost Cloud, whi
repo.
```

$0.80

$637.23

Infracost ★ 10806

Search ⌘ K

Products ∨    Pricing    Resources ∨    About ∨    Sign up / Log in

🏠 > Get started

# Get started

Infracost enables a shift-left approach for cloud costs by providing cost estimates for Terraform **before** deployment. Additionally, it can check for FinOps best practices in accordance with the Well-Architected Frameworks of cloud vendors, and your company's required tag keys/values. This not only saves your team money but also streamlines discussions about costs within the engineering workflow rather than it being a post-deployment consideration. Infracost works with AWS, Azure and Google.

## 1. Install Infracost

Get the latest Infracost release:

**macOS**     macOS/Linux     Windows     Windows     Docker
**brew**      manual          chocolatey  manual

```
brew install infracost

infracost --version # Should show 0.10.38
```

To upgrade Infracost, run `brew update` then `brew upgrade infracost`

AGENT TOOLS
SA Q&A
Generate IaC
Cost Estimation

ORCHESTRATION
Agents

AMAZON BEDROCK
Models to choose
Model Evaluation

CUSTOMIZATION
Prompt Engineering
Knowledge Base (RAG)

INTRO
Use Case

# Bedrock Agent

**Agents for Amazon Bedrock**

Accelerate delivery of generative AI applications

## Create an agent

Use the Bedrock console or SDK to create an agent and provide a description

*"You are a Solutions Architect assistant designed to help customers to design workloads on AWS:*

## Add action groups

Upload API schema so the agent can perform actions (call APIs)

*DescribeDiagram
GenerateTerraform
EstimateCost*

## Add data sources

Configure data sources so the agent can lookup information

*ServiceDocs
ReferenceDiagrams
Whitepapers*

## Interact with the agent

Use natural language to tell the agent to perform a task

*"Estimate monthly cost of running workload"*

# Each Action Group has 3 key elements

## Action Group Description

Overview of actions provided – helps agent know when this action group is relevant

## API Schema

- Rich definition of each action
- Operation name, input parameters, data types, response details
- Helps agents know when to use it, how to call it, and how to use results
- Language agnostic API definition using industry-standard schema

## Lambda Function

- Implementation of each action
- Contains either business logic or wraps microservices, databases, or tools
- Serverless, scalable, secure
- Choice of programming language (Python, C#, JavaScript, Java, …)

AGENT TOOLS
SA Q&A
Generate IaC
Cost Estimation

WRAP-UP
Takeaways
Resources

AMAZON BEDROCK
Models to choose
Model Evaluation

ORCHESTRATION
Agents

CUSTOMIZATION
Prompt Engineering
Knowledge Base (RAG)

INTRO
Use Case

# 3 Key Takeaways

Evaluate different models on cost, speed, and efficiency using Amazon Bedrock's single-API access.

Secure AI Integration with customizable FMs that can be tailored to specific business needs using proprietary data.

Automate complex, multi-step tasks by breaking them down into smaller, manageable actions with Bedrock Agents

# Additional resources

**GitHub repo: Amazon Bedrock samples**

This repository contains pre-built examples to help customers get started with the Amazon Bedrock service including: Knowledge Bases, RAG, Agents, Bedrock Fine-tuning, Security and Governance

**GitHub repo: Build AWS SA using Amazon Bedrock agents**

This repository contains instructions and code samples to help customers This repository contains instructions and code samples for building an AWS Solutions Architect Agent with Amazon Bedrock (SA Q&A, Generate IaC, Estimate Cost)

**PartyRock app: Generative AI Agents For Amazon Bedrock**

This prototype app lets you describe an API, and generate code for Agents for Amazon Bedrock. The generated code includes: an OpenAPI schema providing a rich description of the API, Python Lambda function implementation based on the generated API schema, test suite for the API, delivered as a json array of Lambda event payloads based on the generated API schema.

# Thank you!

**Please complete the session survey in the mobile app**

**Viktoria Semaan**

in linkedin.com/in/semaan

**Viktor Vedmich**

in linkedin.com/in/vedmich