aws marketplace

GENERATIVE AI

# Key prompt engineering strategies to balance cost, performance, and accuracy

rackspace
technology®

# Today's speakers

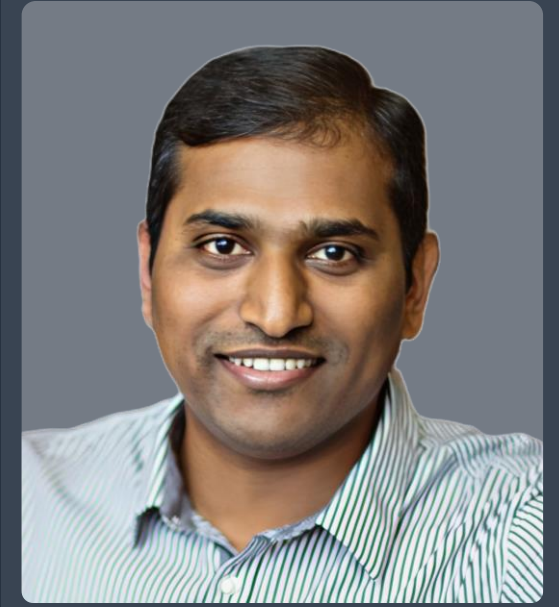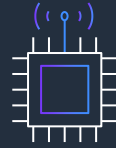| Jacob Newton-Gladstein | Victor Rojo | Pooja Singh | Nirmal Ranganathan |
|---|---|---|---|
| Global Field Deployment Lead, Generative AI Center of Excellence, AWS | Principal Tech Lead, Conversational AI and Generative AI, AWS | Senior Data Science Architect, Rackspace Technology | VP of Engineering, AI Rackspace Technology |

# Innovation can transform industries

GENERATIVE AI

# Generative AI with AWS

## Helping businesses innovate

**Easiest way to build and scale generative AI applications with security and privacy built in**

**Data as your differentiator**

**Generative AI-powered applications to transform user experiences**

**Most performant, low cost infrastructure for generative AI**

# What is prompt engineering?

## What is 10+10?

10+10 = 20

Add guidance

## What is 10+10?

This is an addition problem and the answer to this problem is 20.
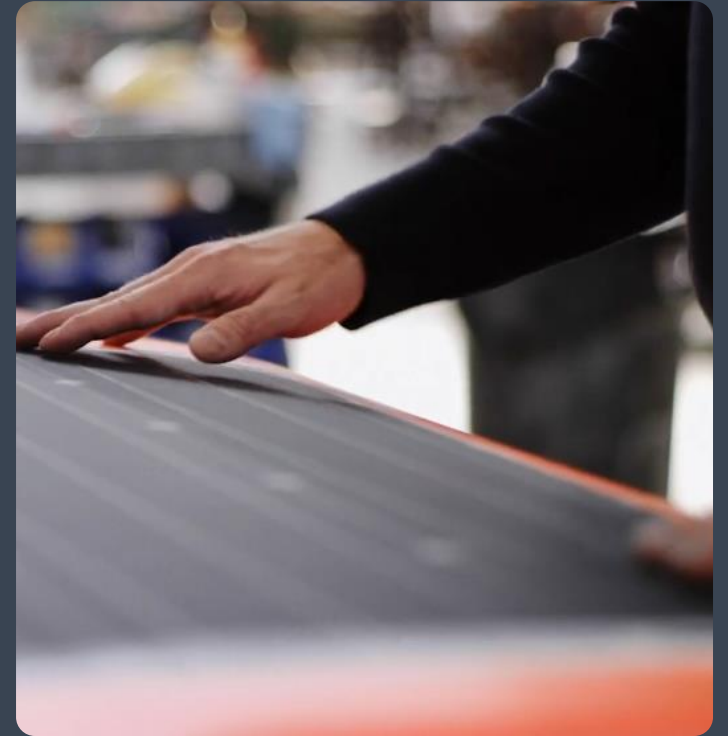
# Prompt engineering benefits



Increase developer control



Improve user experience



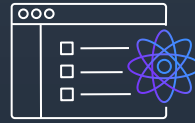Deliver flexibility

# A strategic framework for prompt engineering

**Unambiguous prompts**

**Adequate context within the prompt**

**Balance between targeted information and desired output**

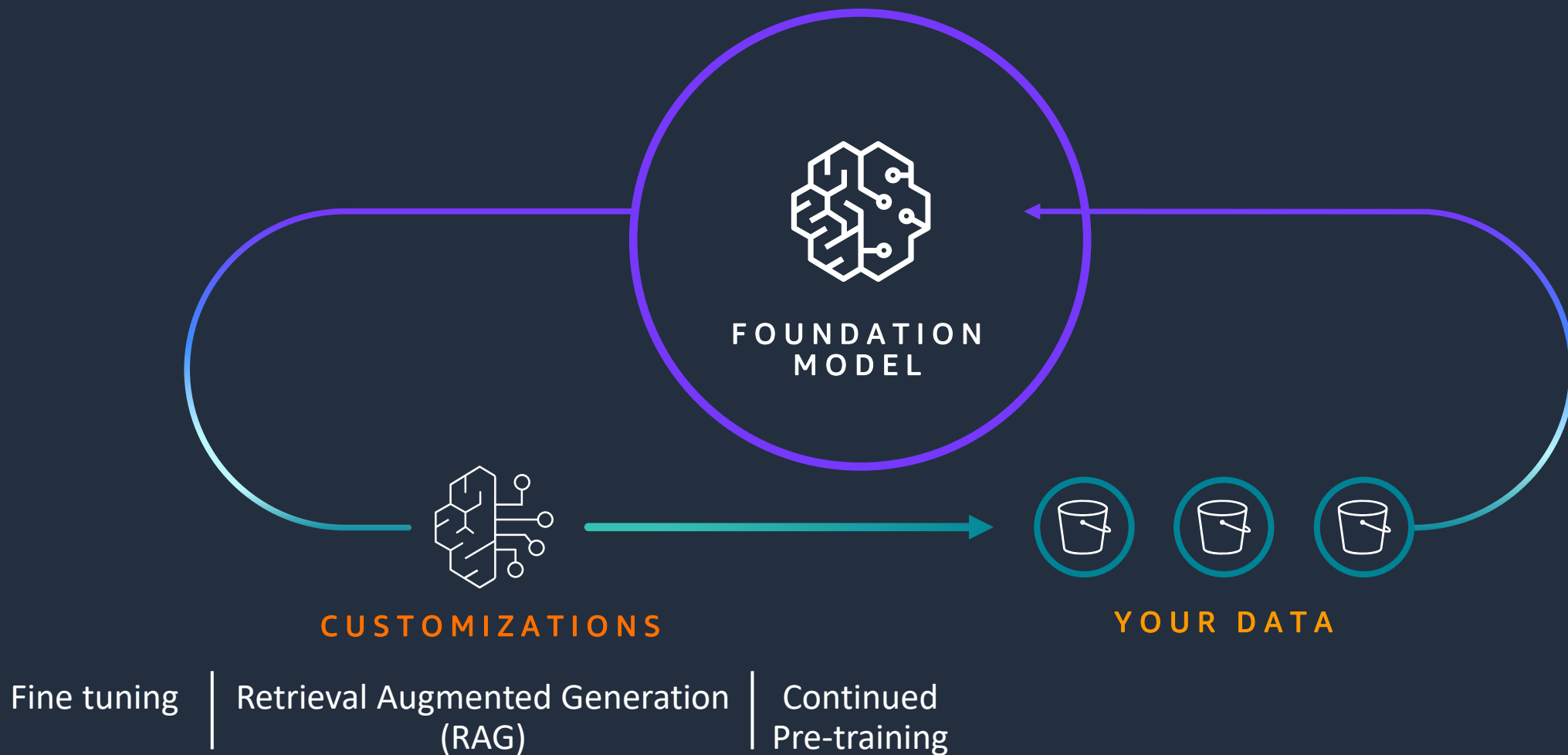**Experiment and refine the prompt**

# Generative AI stack

**APPLICATIONS THAT LEVERAGE LLMs AND OTHER FMs**

**TOOLS TO BUILD WITH LLMs AND OTHER FMs**

**INFRASTRUCTURE FOR FM TRAINING AND INFERENCE**

**FOUNDATION MODEL**

**CUSTOMIZATIONS**

**YOUR DATA**

| Fine tuning | Retrieval Augmented Generation (RAG) | Continued Pre-training |

# Driving outcomes with FAIR™

**rackspace** technology®

# FAIR™ — Foundry for AI by Rackspace

Foundry for AI by Rackspace (FAIR™) is a global incubator focused on **accelerating business value through Responsible AI solutions** across all industries.

- Services
- Capabilities
- Technology | Frameworks
- Partnerships



FAIR™ Industrializing Responsible AI

45+ wins Globally (Americas, EMEA, APAC)

7+ Industries — Cross-Industry
12+ Domains — Cross-Functional
40+ Customers — All maturity levels
Comm, MM, Ent — Cross Segment
5500+ AI Ready
150+ AI Experts — Certified Rackers

FAIR™ Services Portfolio
- Ideate
- Incubate
- Industrialize

AWS Superpowers
aws PARTNER Premier Tier Services
18 AWS Competencies
GenAI Launch Partner

Responsible AI Operating Model Frameworks

Solution Accelerators

Partner Ecosystem
aws · NVIDIA · securiti · Weights & Biases · AIBLE · run:ai

# Responsible AI for a secure and sustainable future

SAFETY,
EXPLAINABILITY, ACCOUNTABILITY,
TRANSPARENCY

## Symbiotic
AI should safely co-exist with us, augmenting our intelligence and making us better at our jobs.

## Sustainable
AI should help us make better decisions while ensuring that its benefits are accessible to everyone.

**RESPONSIBLE AI**

## Secure
AI must be secure, robust and resilient while protecting confidentiality and preventing misuse.

FAIRNESS, ELIMINATING BIAS, ENVIRONMENTAL IMPACT, EQUITY

SECURITY, PRIVACY, RESILIENCE, GOVERNANCE, RISK, COMPLIANCE

# Industrializing AI requires a cloud native execution model

**AI GOVERNANCE**

Cross functional governance team setup to enable agility and manage enterprise risk from AI

## Culture
- Executive Sponsorship
- Business Ownership
- AI as a Co-worker

## People
- Modern hybrid work
- Availability of Cloud, Data and AI Skills
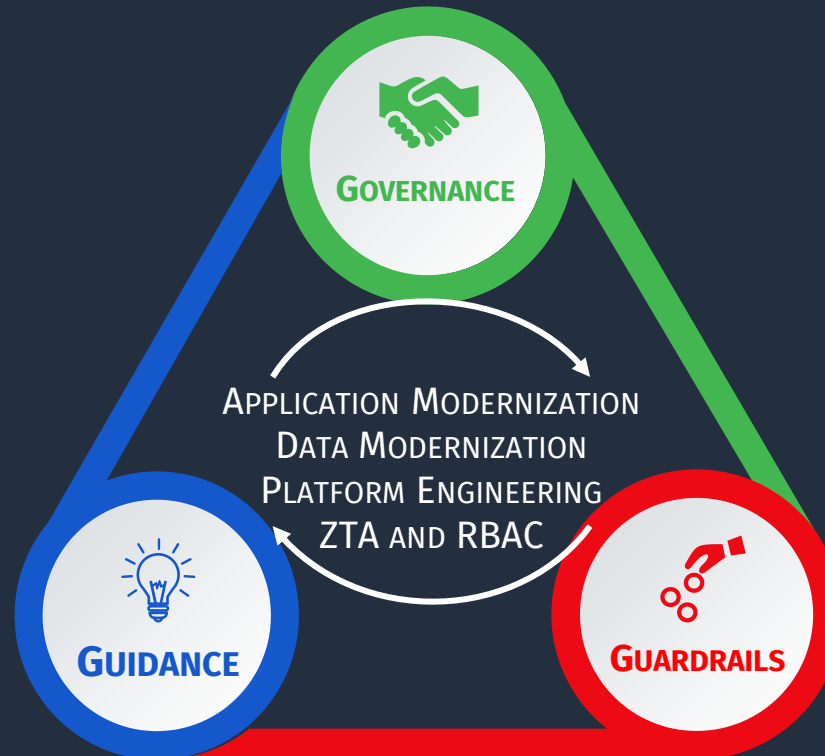- AI Literacy

## Process
- Selecting AI for Responsible Use
- AI Product Management Lifecycle
- AI Service Management

## Technology
- Secure Platform, Identity & Role-based Access Control
- Health and Quality of Supply Chains of AI Data
- Proactive policies to address AI Model Debt

**GOVERNANCE**

APPLICATION MODERNIZATION
DATA MODERNIZATION
PLATFORM ENGINEERING
ZTA AND RBAC

**GUIDANCE**

**GUARDRAILS**

**AI GUIDANCE**

Architectural principles and decision checklists setup to aid decentralized decision making

**AI GUARDRAILS**

Controls and automation setup to prevent accidental and intentional misuse

# FAIR™ Industrialize: Building AI applications and capabilities with a product approach

## Strategy and Innovation
Align AI with the business objectives accounting for growth, risk and sustainability

## Secure AI Platform
Ensure a robust, scalable, and secure AI platform to build next-generation applications

## Supply Chain of Data
Provide consistent, quality data to train and infer the right data-driven decisions and information

## AI Product Development
Re-think design and development of applications infusing user experience and features with AI

## AI Operations
Monitor, manage and operate AI business solutions at scale reliably and responsibly



AI Applications

Strategy & Innovation
Policy Mgmt. · Innovation Mgmt. · AI Literacy · Ethics & Compliance · Sustainability · Hybrid Work · Accountability · Explainability · Transparency

Secure AI Platform
Platform Choice · Data Protection · Guardrails · Threat Management · Access Control · ZTA AI Platform · GRC · Model Selection · Model Security · Model Hub

Supply Chain of Data
Data Sharing · Data Mgmt. · Data Governance · Feature Engineering · Data Provenance · Data Integration · Data Lifecycle · Data Quality · Data Classification

AI Product Development
App Lifecycle · Experience Design · Design Patterns · Model Development · Semantic Programming · Context Mgmt. · Product Decision · AI Safety · Tuning

AI Operations
DevSecOps · DataOps · MLOps · LLMOps · Change Management · Instrumentation · Service Delivery · Incident Mgmt. · FinOps

FAIR™ AI Target Operating Model

# Real-world example of generative AI in action

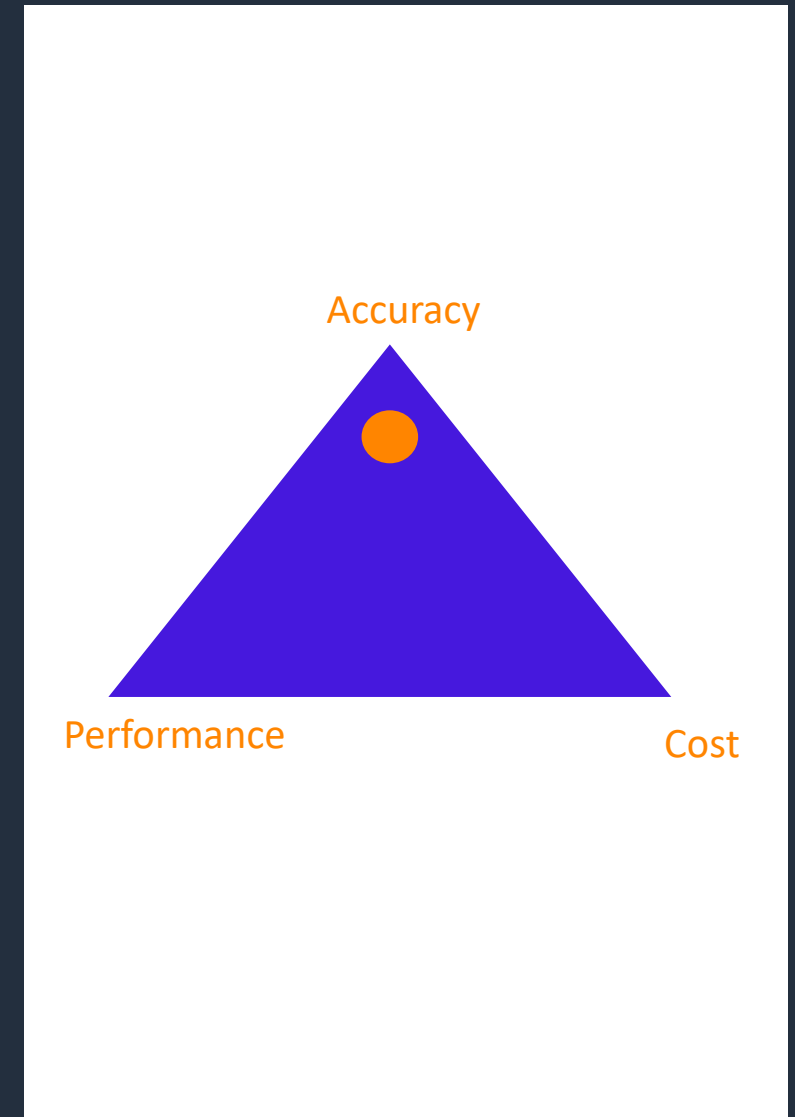## Challenge

- High volume of policy inquiries overwhelming HR and IT teams
- Reduced capacity for high-value tasks
- Inconsistent and time-consuming manual responses
- Need for division-specific policy information access
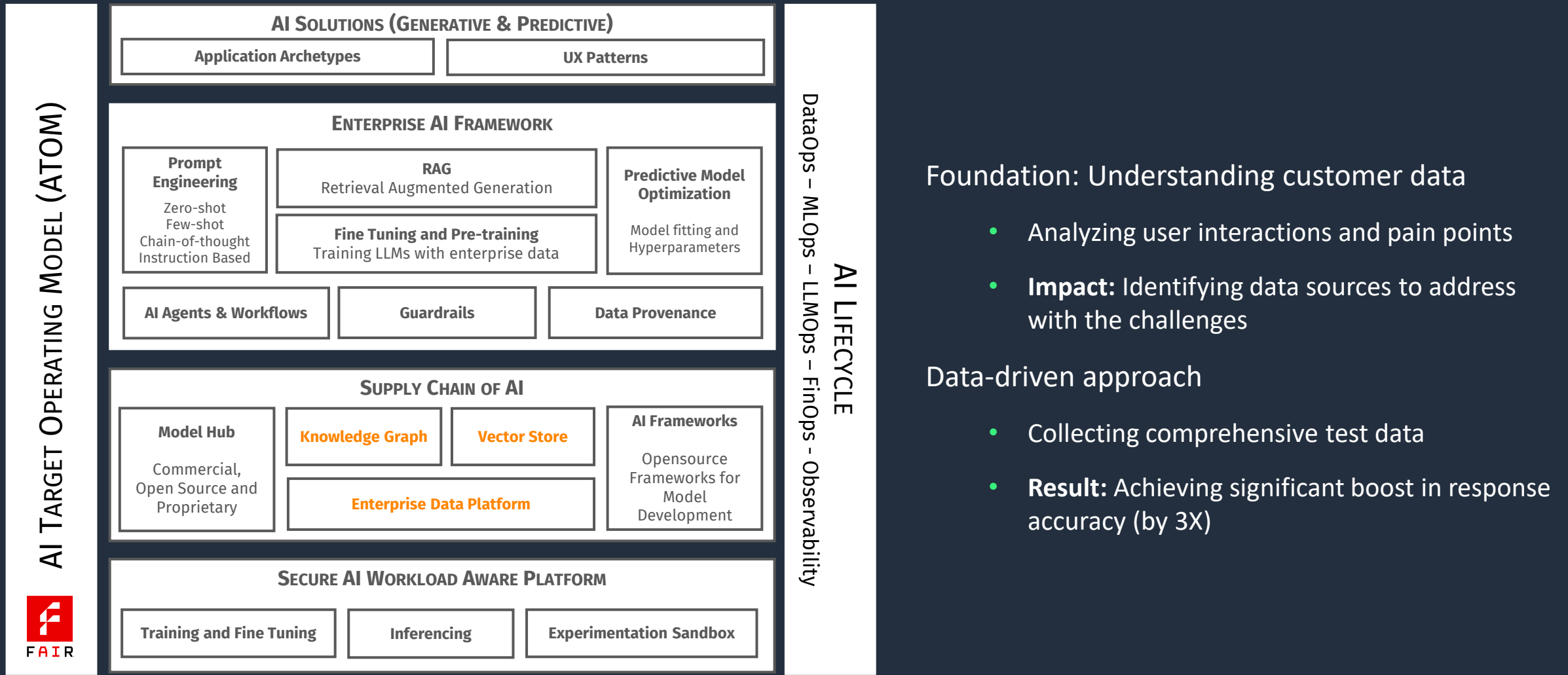
## Solution

- Generative AI-powered RAG chatbot using AWS Bedrock and Claude v3 Sonnet LLM
- Persona-based approach with division-specific metadata tagging
- Scalable architecture supporting 10,000 users
- Robust knowledge base using Amazon S3 and OpenSearch Serverless
- Sophisticated RAG pipeline for accurate and relevant responses

## Success Criteria

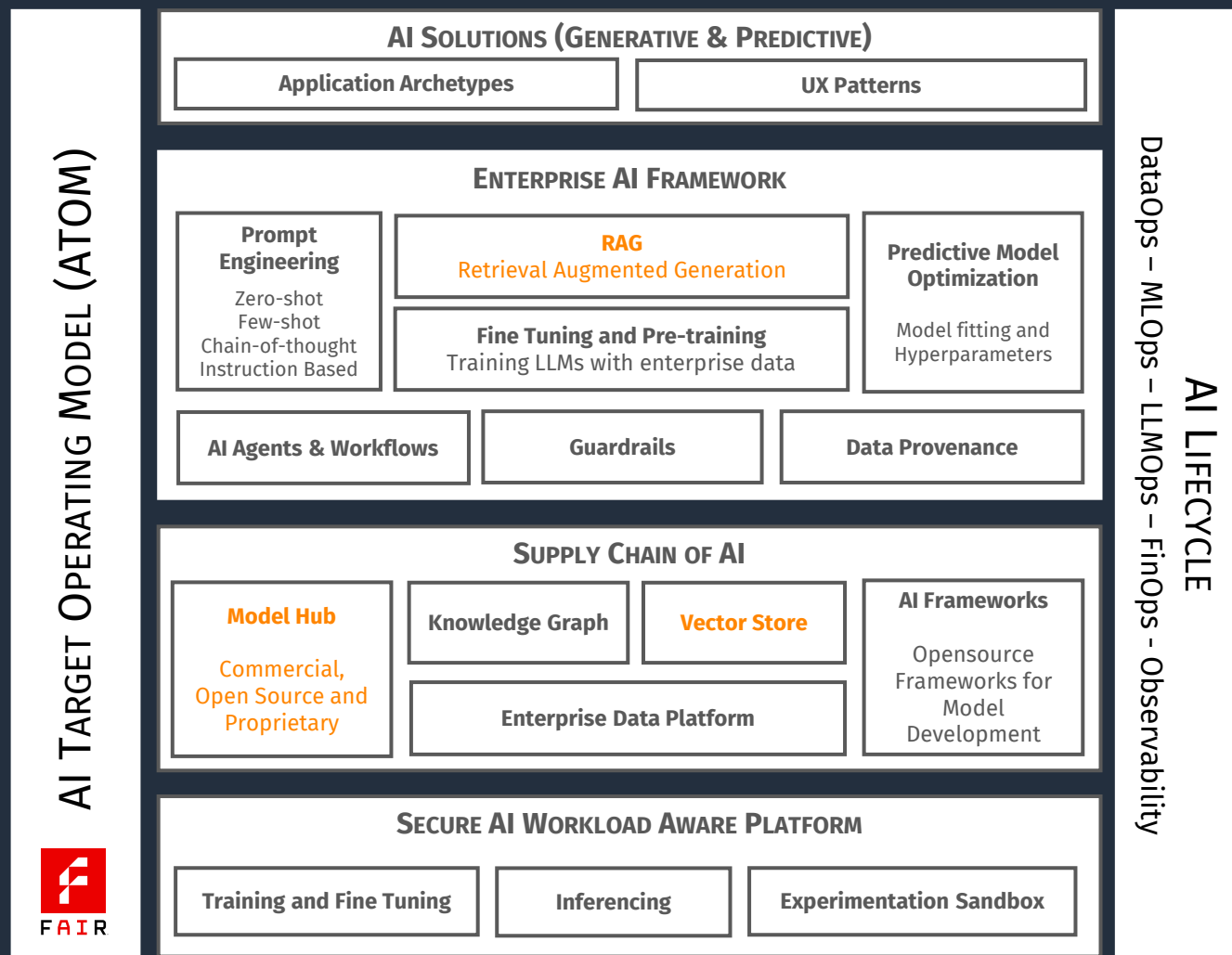- 93% accuracy in initial testing (86 out of 92 test questions correct)
- Projected daily savings of 45 man-hours (30 HR, 15 IT)
- Immediate access to accurate, division-specific policy information
- Improved policy adherence and compliance
- Enhanced overall employee productivity
- Thousands of labor hours redirected to revenue-generating activities

Accuracy

Performance                    Cost

# Supply chain of AI



**AI TARGET OPERATING MODEL (ATOM)**

**AI SOLUTIONS (GENERATIVE & PREDICTIVE)**

| Application Archetypes | UX Patterns |

**ENTERPRISE AI FRAMEWORK**

**Prompt Engineering**
Zero-shot
Few-shot
Chain-of-thought
Instruction Based

**RAG**
Retrieval Augmented Generation

**Fine Tuning and Pre-training**
Training LLMs with enterprise data

**Predictive Model Optimization**
Model fitting and Hyperparameters

| AI Agents & Workflows | Guardrails | Data Provenance |

**SUPPLY CHAIN OF AI**

**Model Hub**
Commercial, Open Source and Proprietary

**Knowledge Graph**

**Vector Store**

**Enterprise Data Platform**

**AI Frameworks**
Opensource Frameworks for Model Development

**SECURE AI WORKLOAD AWARE PLATFORM**

| Training and Fine Tuning | Inferencing | Experimentation Sandbox |

**AI LIFECYCLE**
DataOps – MLOps – LLMOps – FinOps - Observability

## Foundation: Understanding customer data

- Analyzing user interactions and pain points
- **Impact:** Identifying data sources to address with the challenges

## Data-driven approach

- Collecting comprehensive test data
- **Result:** Achieving significant boost in response accuracy (by 3X)

# Optimizing RAG chatbot components



## AI TARGET OPERATING MODEL (ATOM)

### AI SOLUTIONS (GENERATIVE & PREDICTIVE)

| Application Archetypes | UX Patterns |

### ENTERPRISE AI FRAMEWORK

**Prompt Engineering**
Zero-shot
Few-shot
Chain-of-thought
Instruction Based

**RAG**
Retrieval Augmented Generation

**Fine Tuning and Pre-training**
Training LLMs with enterprise data

**Predictive Model Optimization**
Model fitting and Hyperparameters

| AI Agents & Workflows | Guardrails | Data Provenance |

### SUPPLY CHAIN OF AI

**Model Hub**
Commercial, Open Source and Proprietary

**Knowledge Graph**

**Vector Store**

**AI Frameworks**
Opensource Frameworks for Model Development

**Enterprise Data Platform**

### SECURE AI WORKLOAD AWARE PLATFORM

| Training and Fine Tuning | Inferencing | Experimentation Sandbox |

**AI LIFECYCLE**
DataOps – MLOps – LLMOps – FinOps - Observability

## Retrieval tuning

- Embedding model selection with efficient vector store (HNSW algorithm)

- Advanced reranking system (such as Cohere ReRanker) to prioritize documents

- **Outcome:** Accelerated retrieval and precision gains

## Generation tuning

- Iterative and advanced prompt tuning techniques

- LLM selection and context-aware answer generation

- **Impact:** Enhanced answer relevance scores

# Boosting accuracy with prompt engineering

## AI TARGET OPERATING MODEL (ATOM)

### AI SOLUTIONS (GENERATIVE & PREDICTIVE)

| Application Archetypes | UX Patterns |
|---|---|

### ENTERPRISE AI FRAMEWORK

**Prompt Engineering**

Zero-shot
Few-shot
Chain-of-thought
Instruction Based

**RAG**
Retrieval Augmented Generation

**Fine Tuning and Pre-training**
Training LLMs with enterprise data

**Predictive Model Optimization**

Model fitting and Hyperparameters

| AI Agents & Workflows | Guardrails | Data Provenance |
|---|---|---|

### SUPPLY CHAIN OF AI

**Model Hub**

Commercial, Open Source and Proprietary

**Knowledge Graph**

**Vector Store**

**AI Frameworks**

Opensource Frameworks for Model Development

**Enterprise Data Platform**

### SECURE AI WORKLOAD AWARE PLATFORM

| Training and Fine Tuning | Inferencing | Experimentation Sandbox |
|---|---|---|

**AI LIFECYCLE**
DataOps – MLOps – LLMOps – FinOps - Observability

**FAIR**

## Advanced prompt engineering techniques:

- Instructions for the prompts with guardrails
- Few-shot examples for persona-based responses to the user questions
- Caching mechanism for frequently asked questions to improve response times
- **Result:** Improvement in accuracy to 93% and user satisfaction

## Iterative refinement

- Continuous testing and optimization
- **Outcome:** Consistent accuracy improvements with each iteration

## Accuracy-performance-cost optimization

- Streamlining resource utilization and model efficiency
- **Final result:** Accuracy enhancement while balancing performance and cost

# Prompt engineering examples

```
<instruction>
        You are a policy support agent trained on company policies.
        Use the provided context to answer the question.
        If you don't know, say so. Don't make up answers.
        Give the answer from the user persona if available, otherwise use the corporate persona.
        Only use relevant context chunks for the question and user persona.
        Explain reasoning in <reasoning> tags and give the answer in <answer> tags.
        Don't mention context numbers in the <answer> tags.
</instruction>

<example>
        Context 1: [Policy about topic A1]
        Context 1 Persona: factory-division

        Context 2: [Policy about topic A2]
        Context 2 Persona: corporate-division

        Context 3: [Policy about topic B]
        Context 3 Persona: retail-division

        Question: [Question about policy topic related to A1]
        User Persona: factory-division

        Output:
                <reasoning>
                                [Explanation of why Context 1 is relevant]
                </reasoning>
                <answer>
                                [Answer based on Context 1 for factory-division persona]
                </answer>
</example>
```

```
<example>

        Context 1: [Policy about attendance and notification procedures]

        Context 1 Persona: retail-division

        Context 2: [Policy about weather-related closures]

        Context 2 Persona: retail-division

        Context 3: [Policy about time off requests]

        Context 3 Persona: corporate-division

        Question: [Question about calling in sick for a shift]

        User Persona: retail-division

        Output:

                <reasoning>

                                [Explanation of why Context 1 is most relevant for this question]

                </reasoning>

                <answer>

                                [General answer about notifying manager of absence, based on
        retail persona policy]

                </answer>
</example>
```

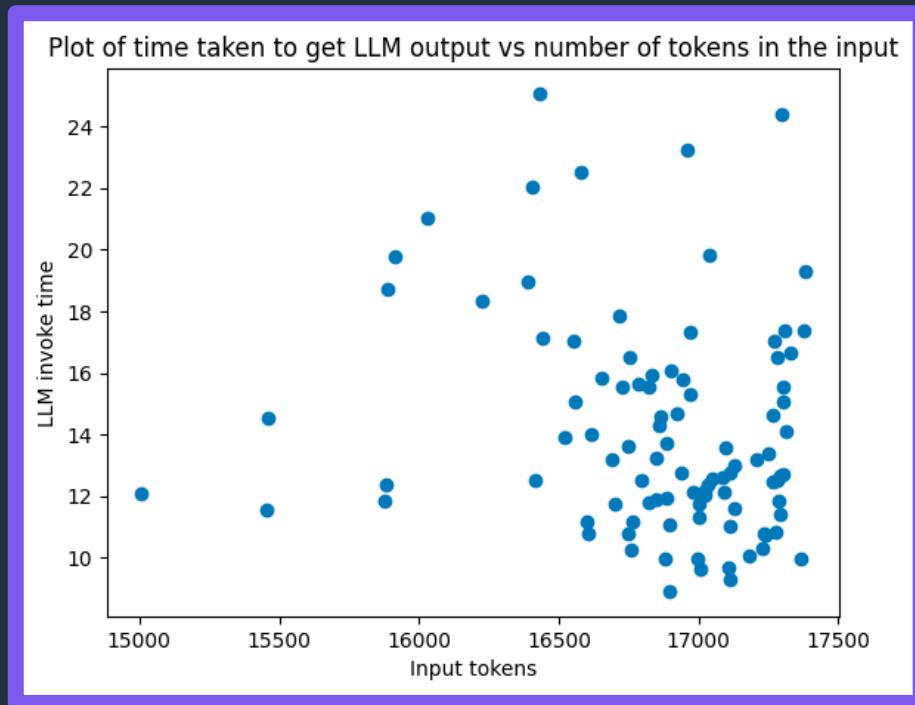# Optimizing vector store retrieval: Finding the ideal K for maximum recall
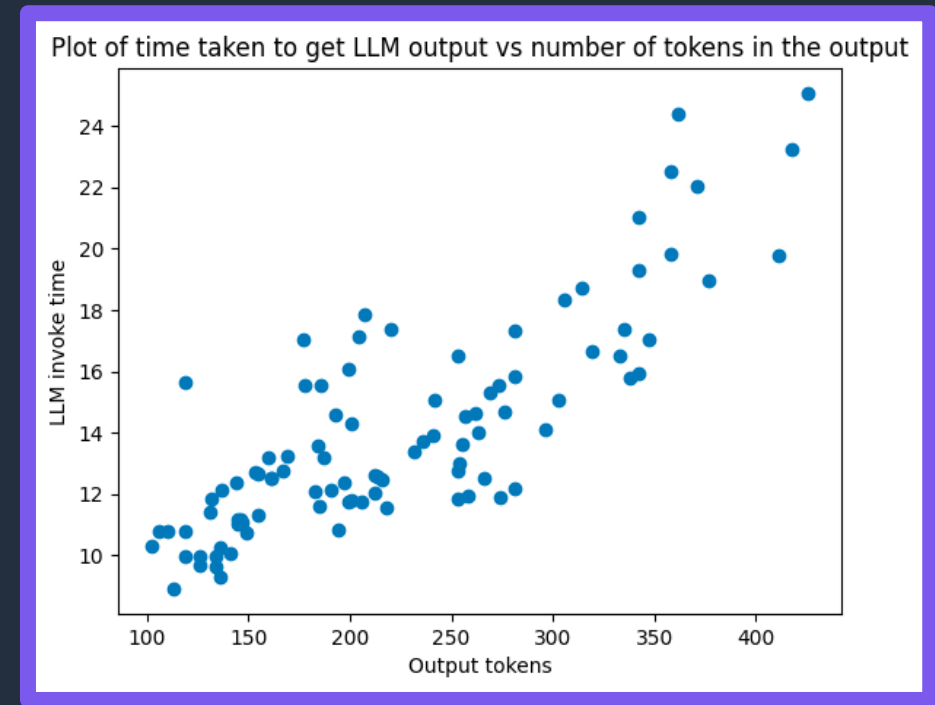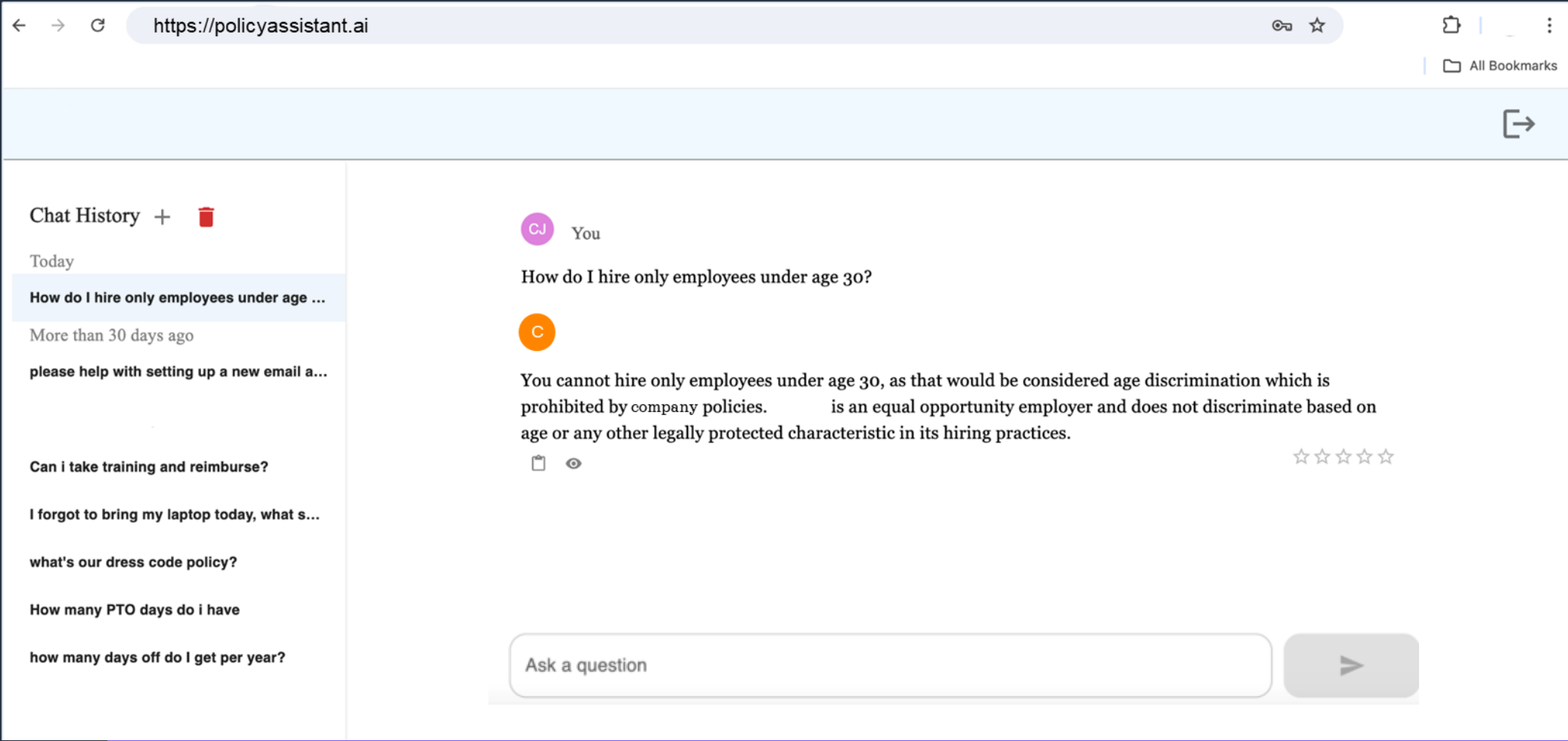


With Cosine Similarity



With Cohere ReRanker

# Optimizing RAG performance: Token count impact on response time



Plot of time taken to get LLM output vs number of tokens in the input



Plot of time taken to get LLM output vs number of tokens in the output

- More input tokens generally increase retrieval time
- Increases cost for the context passed to FMs

- Increasing output tokens extends generation time, increasing overall latency
- Also, higher generation costs

# Policy chat bot

# Lessons learned

Hybrid approach combining few-shot, zero-shot, and role-playing prompts

Enhanced policy compliance across the organization

Accelerated organizational growth through efficient, ethical AI adoption

Persona-based responses tailored to different divisions (Corporate, Retail, Factory)

Ethical AI implementation with robust governance and guardrails

Significant reduction in response time for employee inquiries – 45 hours saved daily!

# Q&A

## Jacob Newton-Gladstein

Global Field Deployment Lead,
Generative AI Center
of Excellence, AWS

## Victor Rojo

Principal Tech Lead,
Conversational AI and
Generative AI, AWS

## Pooja Singh

Senior Data Science Architect,
Rackspace Technology

## Nirmal Ranganathan

VP of Engineering, AI
Rackspace Technology

# How you can engage Rackspace Technology

Get better results throughout your generative AI journey with Rackspace Technology.
See our offering on AWS Marketplace.

**Visit Rackspace on AWS Marketplace**

## Generative AI Ideation | FAIR™

Determine business goals and define a top-priority generative AI use case in three weeks



## Generative AI Incubate for Amazon Q

Establish the feasibility of AI in your business with an Amazon Q MVP for your first AI use case

# What is AWS Marketplace?

**aws marketplace**

- Over **20K** listings
- **4K+** ISVs (free, BYOL, or commercial)
- Deployed in **31** regions
- **300K+** monthly active customers
- Over **2.5M** current subscriptions
- Offers **70+** product categories

- Deploy software on demand
- Flexible consumption and contract models
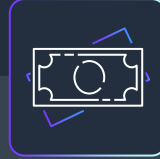- Easy and secure deployment, almost instantly
- Simplified billing



**aws marketplace**

# How can you get started?

### Find
A breadth of generative
AI solutions

Applications

Foundation models

Tooling

Professional services

Data sets

### Buy
Through flexible
pricing options

Free trial

Pay-as-you-go

Budget alignment

Private offers

Billing consolidation

Enterprise Discount Program

Private Marketplace

### Deploy
With multiple
deployment options

AI/ML models

Amazon Machine Image

Containers

CloudFormation template

Amazon EKS/Amazon ECS

SaaS

AWS Data Exchange

aws marketplace

# Upcoming spotlight series
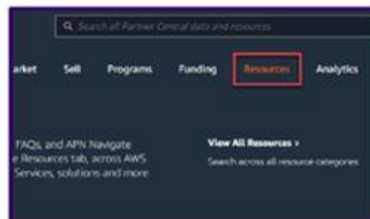
# The Generative AI Center of Excellence
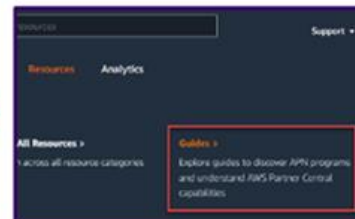
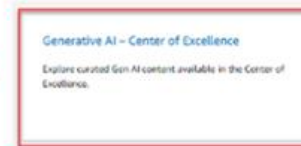190+ Generative AI Assets

80,000+ Interactions

1300+ AWS Partners



Partner Central > Resources

Guides

Generative AI – Center of Excellence

Partner Central CoE Page

# Thank you!