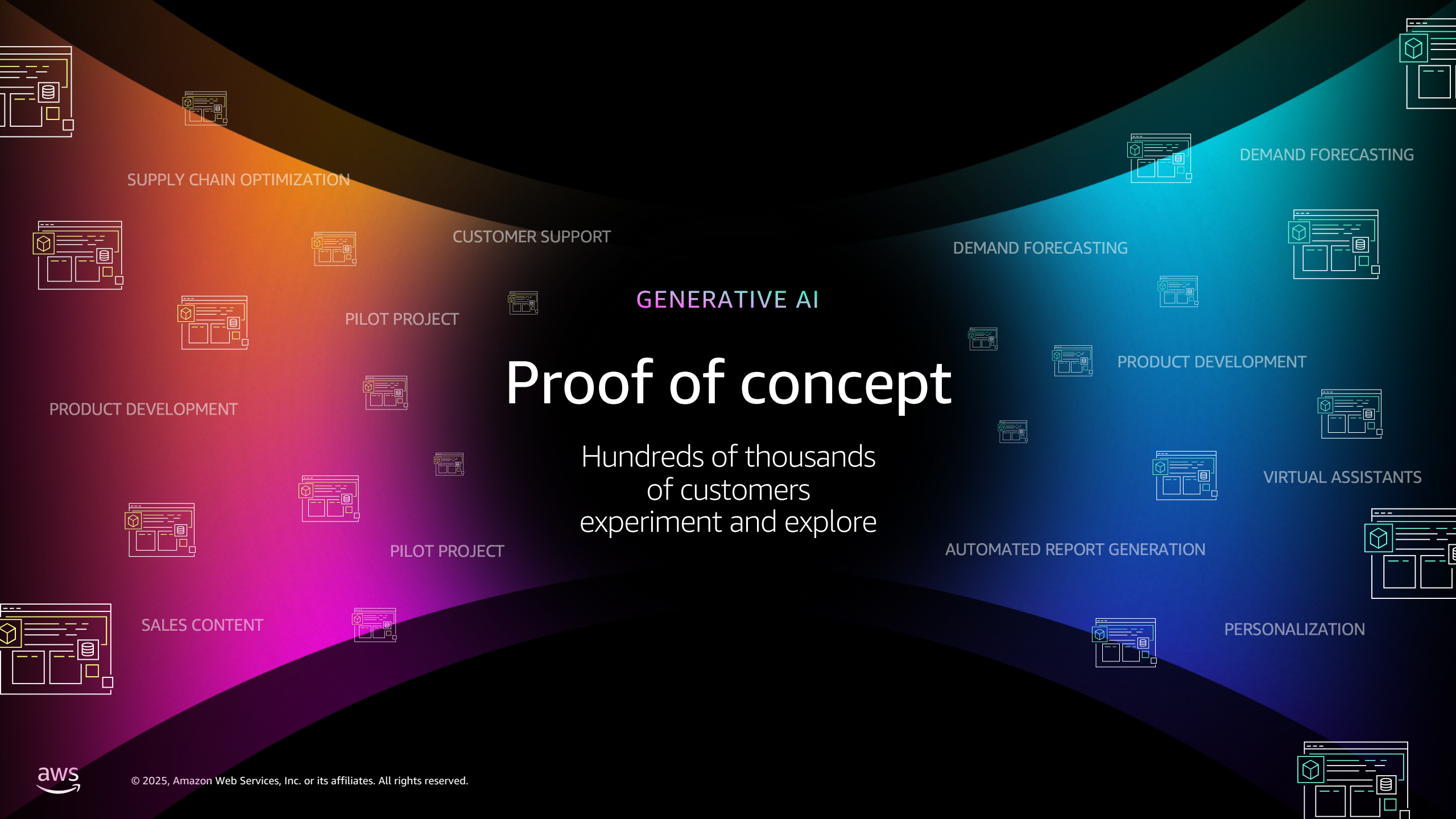


# Generative AI in Action: From prototype to production

**Dragoș Mădărășan**

Solutions Architecture Team Lead,  
AWS





SUPPLY CHAIN OPTIMIZATION

DEMAND FORECASTING

CUSTOMER SUPPORT

DEMAND FORECASTING

GENERATIVE AI

PILOT PROJECT

PRODUCT DEVELOPMENT

# Proof of concept

Hundreds of thousands  
of customers  
experiment and explore

PRODUCT DEVELOPMENT

VIRTUAL ASSISTANTS

PILOT PROJECT

AUTOMATED REPORT GENERATION

SALES CONTENT

PERSONALIZATION



ACT 1

Proliferation of  
prototypes

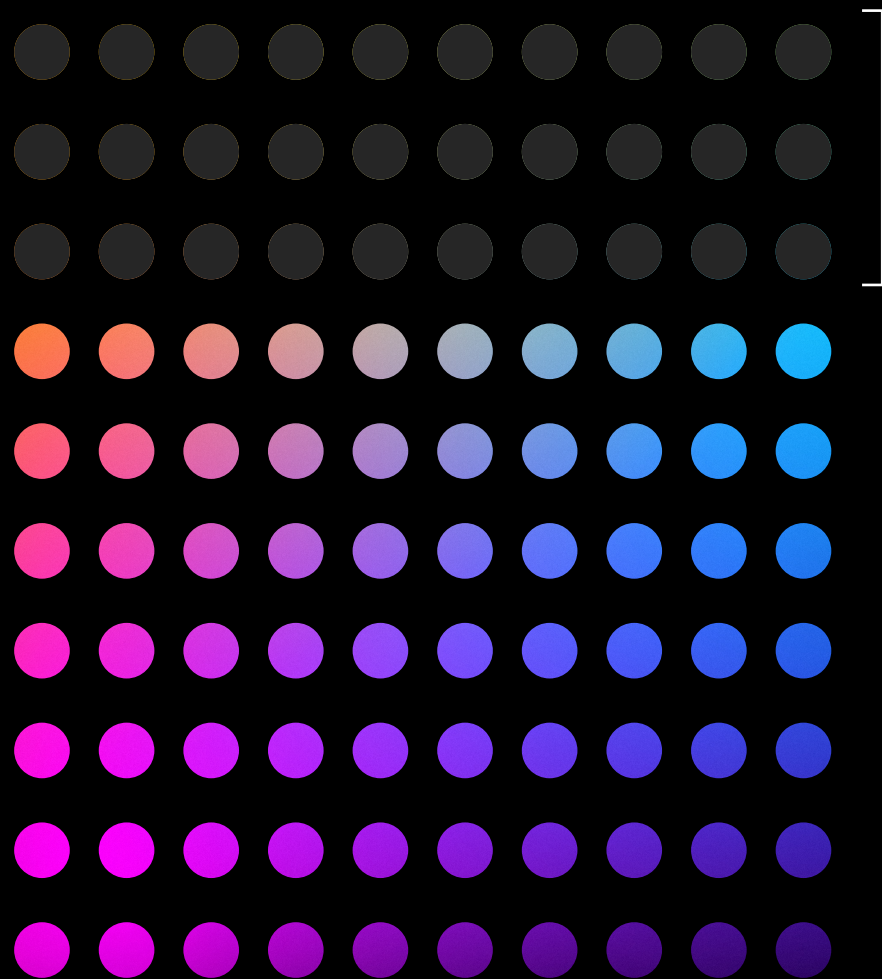
Generative AI

ACT 1

**Proliferation of  
prototypes**

ACT 2

**Acceleration of  
production**



Gartner® predicts

# 30%

of generative AI projects will  
be abandoned after proof of  
concept by the end of 2025

# Four key drivers for success





# Models

SUMMARIZATION



FAST



LOW COST



REASONING



# No one model to rule them all

CODE  
INTERPRETATION



INTELLIGENT



ANALYSIS





# Amazon SageMaker AI

Build, train, deploy machine learning models for any use case with fully managed infrastructure, tools, and workflows

NEW

# Amazon Nova

State-of-the-art foundation models that deliver frontier intelligence and industry leading price performance





# Amazon Bedrock

The easiest way to build and scale  
generative AI applications

# Amazon Bedrock

## Broadest selection of models

### CUSTOM MODEL IMPORT

Leverage your customized models on Amazon Bedrock

**AI21 labs**

Effective reasoning & rapid analysis for long context windows

**JAMBA**

**amazon**

Frontier multimodal intelligence at low-latency, Agent & RAG Applications, high-quality image & video generation

**AMAZON NOVA**

**ANTHROPIC**

Advanced reasoning & coding capabilities, including computer use skills

**CLAUDE**

**cohere**

Multimodal search & advanced retrieval powering multilingual knowledge agents

**COMMAND  
EMBED  
RERANK**

**Luma**

High-quality video generation from text & images

**LUMA RAY 2**

**Meta**

Advanced image & language reasoning

**LLAMA**

**MISTRAL  
AI**

Knowledge summarization expert agents, & code completion

**MISTRAL  
MIXTRAL**

**poolside**

Software engineering AI for large enterprises

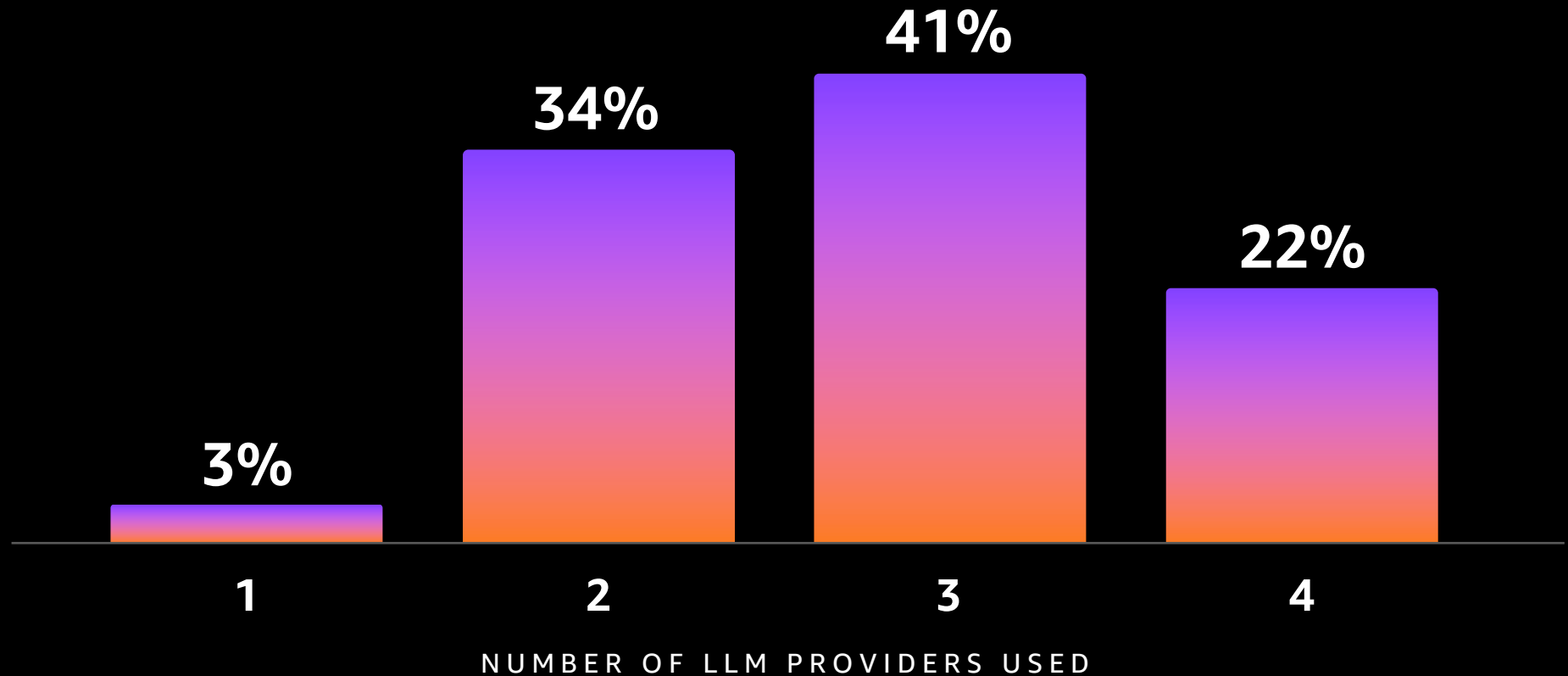
**MALIBU  
POINT**

**stability.ai**

High-quality AI image generation, easily deployable at scale

**STABLE  
DIFFUSION  
STABLE  
IMAGE**

# Enterprises are deploying models from multiple model providers



NEW

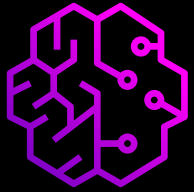
# Amazon Bedrock Marketplace

Discover, test, and use over 100 popular, emerging, and specialized models in Amazon Bedrock





# Cost



# Amazon SageMaker HyperPod

**Reduce time to train foundation models** by up to 40% and scale across more than a thousand AI accelerators efficiently

# Inference cost challenges



NEW

# Amazon Bedrock supports prompt caching

Cache repetitive context in prompts across multiple API calls

AVAILABLE IN PREVIEW

NEW

# Amazon Bedrock Intelligent Prompt Routing

Automatically route prompts to models to optimize response quality and lower costs

AVAILABLE IN PREVIEW



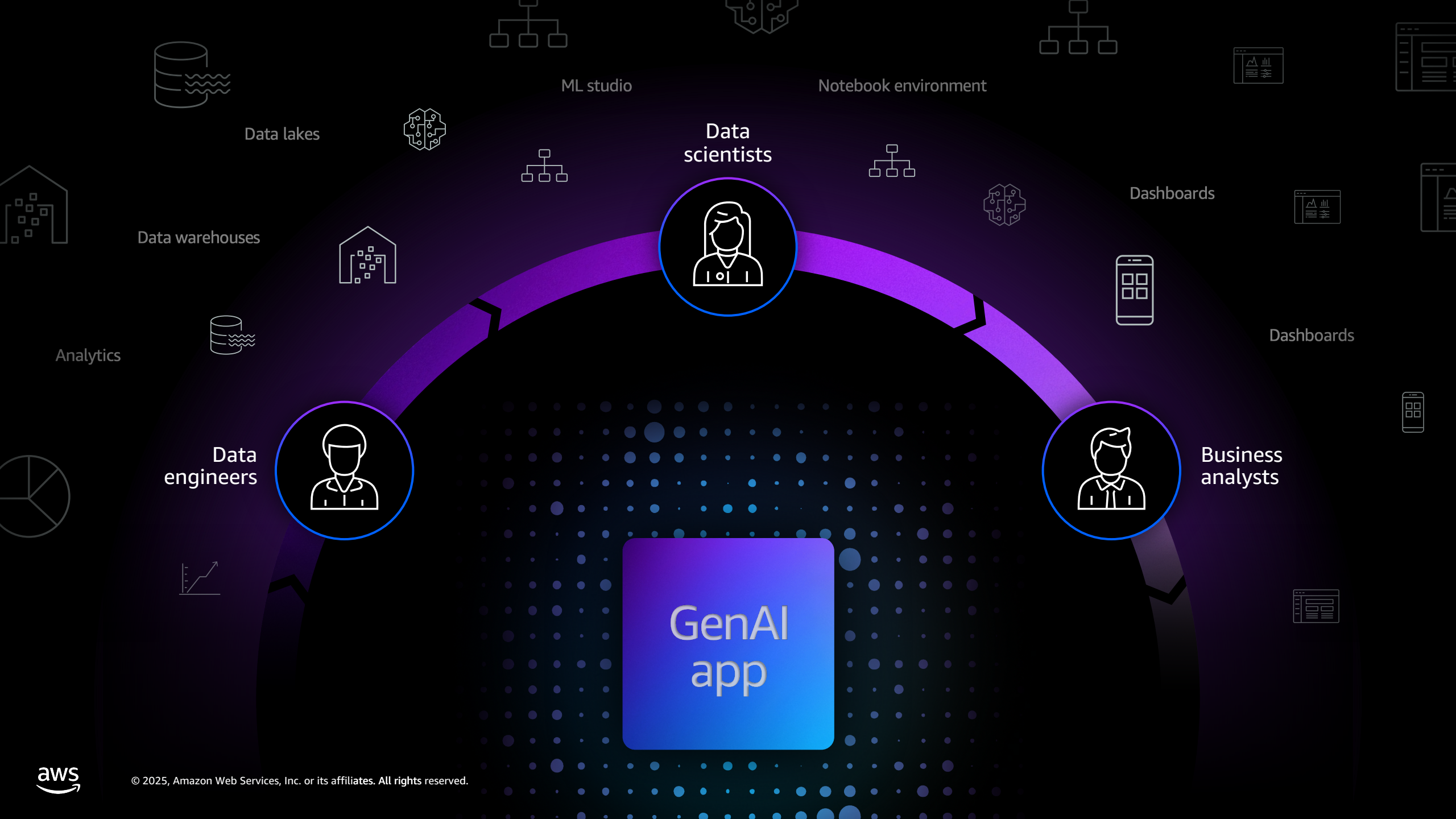


# Data



# Amazon Bedrock Knowledge Bases

Fully managed support for end-to-end  
RAG workflow

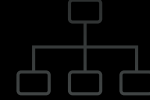


Data lakes



ML studio

Notebook environment



Dashboards



Dashboards



Data warehouses



Analytics



Data engineers



Data scientists



Business analysts



GenAI app





# Trust

# ISO 42001



Amazon  
Bedrock



Amazon Q  
for Business



Amazon  
Textract



Amazon  
Transcribe

Accredited Certification

Up to

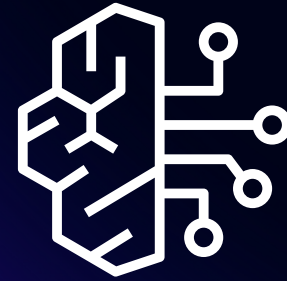
85%

More harmful content blocked  
with our safety features

---

75%

Hallucinated responses filtered



Amazon Bedrock  
Guardrails

NEW

# Amazon Bedrock now supports multi-agent collaboration

Build powerful multi-agent systems—  
faster, easier, and more cost-effectively

AVAILABLE IN PREVIEW



# Four key drivers for Act 2

Trust



Data



Cost



Models





**The best place** to build your  
generative AI applications



# Thank you!

**Dragoș Mădărășan**

